

# Parameters or Privacy: A Provable Tradeoff Between Overparameterization and Membership Inference

Jasper Tan, Blake Mason, Hamid Javadi, Richard G. Baraniuk  
Rice University

## Abstract

A surprising phenomenon in modern machine learning is the ability of a highly *overparameterized* model to generalize well (small error on the test data) even when it is trained to memorize the training data (zero error on the training data). This has led to an arms race towards increasingly overparameterized models (c.f., deep learning). In this paper, we study an underexplored hidden cost of overparameterization: the fact that overparameterized models may be more vulnerable to *privacy attacks*, in particular the *membership inference* attack that predicts the (potentially sensitive) examples used to train a model. We significantly extend the relatively few empirical results on this problem by theoretically proving for an overparameterized linear regression model in the Gaussian data setting that membership inference vulnerability increases with the number of parameters. Moreover, a range of empirical studies indicates that more complex, nonlinear models exhibit the same behavior. Finally, we extend our analysis towards ridge-regularized linear regression and show in the Gaussian data setting that increased regularization also increases membership inference vulnerability in the overparameterized regime.

## 1 Introduction

As more machine learning models are being trained on sensitive user data (e.g., customer behavior, medical records, personal media), there is a growing concern that these models may serve as a gateway for malicious adversaries to access the models’ private training data [1, 2]. For example, the possibility of performing *membership inference* (MI), the task of identifying whether a specific data point was included in a model’s training set or not, can be greatly detrimental to user privacy. In settings like healthcare, knowledge of mere inclusion in a training dataset (e.g., hospital visitation records) already reveals sensitive information. Moreover, MI can enable the extraction of verbatim training data from a released model [1].

The general motivation for this paper is the recent realization that privacy issues around machine learning might be exacerbated by today’s trend towards increasingly *overparameterized models* that have more parameters than training data points and so can be trained to memorize (attain zero error on) the training data. Surprisingly, some overparameterized models (e.g., large regression models [3], massive deep networks [4, 5]) generalize extremely well [6, 7]. Limited empirical evidence suggests that overparameterization may lead to greater privacy vulnerabilities [1, 8–11]. However, there has been little to no analytical work on this important problem.

*In this paper, we take the first steps towards an analytical understanding of how the number of parameters in a machine learning model relates to MI vulnerability.* In a theoretical direction, we prove for linear regression on Gaussian data in the overparameterized regime that increasing the number of parameters of the model increases its vulnerability to MI (Theorem 3.2). In a supporting empirical direction, we demonstrate that the same behavior holds for a range of more complex models: a latent space model, a time-series model, and a nonlinear random ReLU features model (Section 5).

We also extend our analysis towards the techniques of ridge regularization and noise addition in Section 4. We first prove that MI vulnerability similarly increases with the number of parameters for ridge-regularized models. Surprisingly, we also demonstrate analytically that additional regularization *increases* this vulnerability in overparameterized linear models. Hence, ridge regularization is not always an effective

defense against MI and can even be harmful. Afterwards, we show that the privacy-utility trade-off induced by reducing the number of parameters of a linear regression model is equivalent to that obtained by adding independent Gaussian noise to the output of the model when using all available parameters.

Overall, our analyses show that there are settings where reducing the number of parameters of a highly overparameterized model is a simple strategy to protect a model against MI. As the trend towards increasingly overparameterized machine learning models accelerates, *our results make the case for less overparameterization when privacy is a concern.*

**Related work.** The problem of membership inference has received considerable attention, and we refer the reader to [12] for a survey of the prior art. Many works have demonstrated how deep neural networks can be highly vulnerable to MI attacks (cf., [13–16] for example). Various types of attacks, such as shadow models [13], confidence-based attacks [17, 18], and label-only attacks [19] have risen to further expose privacy vulnerability of machine learning models. Our analysis in this work is based on optimal MI attacks, a framework also employed by other studies of MI [9, 10, 14].

While we study the effect parameterization has on MI vulnerability, several studies examine how other aspects of the machine learning pipeline, such as data augmentation [20], dropout [15], sparsity [21], and ensembles [22] affect MI vulnerability. [16] show that pruning a neural network can defend against MI, lending further experimental support to the link between overparameterization and membership inference. [9] analyzes the relationship between overfitting (as measured by generalization error) and membership inference but does not connect this to the number of parameters. Much recent work has shown though that overparameterization does not always lead to increased generalization error and indeed sometimes can even decrease it [3, 6, 7, 23], suggesting that more work is needed linking overparameterization and membership inference beyond mere generalization error. Despite some empirical evidence that overparameterized models are vulnerable to MI attacks [1, 11], to the best of our knowledge, there has been no theoretical study of the effect of this connection.

Differential privacy (DP) [24] is another popular framework frequently used to study the privacy properties of machine learning algorithms, and models that have DP also enjoy MI guarantees [14]. [25] bounds the minimum dataset size as a function of dimension to ensure privacy, however its results do not extend to data of continuous values (e.g., regression data) as is done in this work and does not account for the randomness in the data generating process, which as argued in [14], is a necessary condition to bound the optimal MI risk.

While our theoretical analysis focuses on linear regression, many works have studied overparameterized linear models as an interpretable and tractable test case of more complex overparameterized models such as deep networks [3, 6, 23, 26, 27], suggesting links of our results to more complex nonlinear models.

**Contributions.** In this paper, we establish the first theoretical connection between overparameterization and vulnerability to membership inference attacks. Our contributions are as follows:

1. We derive the optimal MI accuracy for linear regression on Gaussian data and show that increasing the number of parameters increases the vulnerability to MI.
2. We empirically show for more complex models that increasing the number of parameters also increases MI vulnerability.
3. We theoretically show for overparameterized linear models that increased ridge regularization increases MI vulnerability in the Gaussian data setting.
4. We use our analysis to show that reducing the number of parameters of a linear regression model yields an equivalent privacy-utility trade-off as adding independent noise to the output of a model using all available parameters.

Our code is provided in [https://github.com/tanjasper/parameters\\_or\\_privacy](https://github.com/tanjasper/parameters_or_privacy).

## 2 The Membership Inference Problem

We now introduce our notation and formally state the membership inference (MI) problem. Let  $D \in \mathbb{Z}^+$  denote the data dimension. Let  $z = (x, y) \in \mathbb{R}^D \times \mathbb{R}$  denote a data point, and let  $\mathcal{D}$  be a distribution

over the data points. Consider a set  $\mathcal{S}$  sampled  $\mathcal{S} \sim \mathcal{D}^n$  of size  $n \in \mathbb{Z}^+$  data points and denote it by  $\mathcal{S} = \{z_1, z_2, \dots, z_n\} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ . Let  $f_S$  represent a machine learning model obtained by applying a training algorithm  $T$  on the dataset  $\mathcal{S}$ . In particular,  $f_S$  is a deterministic function in  $\mathbb{R}^D \rightarrow \mathbb{R}$ . Typically, it is a function that minimizes  $\sum_{i=1}^n \ell(f(x_i), y_i)$  over a family of functions  $f$  for some loss function  $\ell : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$ . We denote  $\hat{y}_i = f_S(x_i)$ .

Given a data point  $z_0$  and a trained model  $f_S$ , *membership inference* is the task of identifying if  $z_0$  was included in the set  $\mathcal{S}$  used to train  $f_S$ . Different MI attacks are characterized by the types of information accorded to the adversary. In this work, we focus on a blind black-box setting where the adversary has access to  $\mathbf{x}_0$ ,  $f_S(\mathbf{x}_0)$ ,  $n$ , the training algorithm  $T$ , and the data distribution  $\mathcal{D}$ , but not to the ground truth  $y_0$ , the loss value  $\ell(f_S(\mathbf{x}_0), y_0)$ , or any learned parameters of  $f_S$ . Thus, the adversary is a function  $A : \mathbb{R}^d \times \mathbb{R} \rightarrow \{0, 1\}$  such that given  $(\mathbf{x}_0, f_S(\mathbf{x}_0))$ , it outputs 0 or 1, representing its prediction as to whether  $\mathbf{x}_0$  is a member of  $\mathcal{S}$  based on  $f_S(\mathbf{x}_0)$ .

There are two main interests in studying blind MI wherein the adversary does not have access to the loss value. First, there are many realistic scenarios where the ground truth variable of interest is unknown for the general population, and hence the need for a model  $f_S$  to predict it. Second, since we study the interpolation regime where the training loss can be driven down to exactly 0, MI becomes trivial given the loss value: an adversary that predicts any data point that achieves 0 loss as a member will have perfect accuracy. Multiple works in the literature also study this blind setting [18, 28]. We emphasize that our results on the blind adversary lower bound the performance of adversaries that additionally have access to more information.

MI is often defined as an experiment to facilitate precise analysis, and we follow the setup of [9] summarized below.

**Experiment 1. Membership inference experiment.** Given a data distribution  $\mathcal{D}$ , an integer  $n$ , a machine learning training algorithm  $T$ , and an adversary  $A$ , the MI experiment is as follows:

1. Sample  $\mathcal{S} \sim \mathcal{D}^n$ .
2. Apply the training algorithm  $T$  on  $\mathcal{S}$  to obtain  $f_S$ .
3. Sample  $m \in \{0, 1\}$  uniformly at random.
4. If  $m = 0$ , sample a data point  $(\mathbf{x}_0, y_0) \sim \mathcal{D}$ . Else, sample a data point  $(\mathbf{x}_0, y_0) \in \mathcal{S}$  uniformly at random.
5. Observe the adversary's prediction  $A(\mathbf{x}_0, f_S(\mathbf{x}_0)) \in \{0, 1\}$ .

We wish to quantify the optimal performance of the adversary  $A$  and turn to the popular *membership advantage* metric defined in [9]. The membership advantage of  $A$  is the difference between the true positive rate and false positive rate of its predictions.

**Definition 2.1** ([9]). The *membership advantage* of an adversary  $A$  is:

$$\text{Adv}(A) := \mathbb{P}(A(\mathbf{x}_0, f_S(\mathbf{x}_0)) = 1 \mid m = 1) - \mathbb{P}(A(\mathbf{x}_0, f_S(\mathbf{x}_0)) = 1 \mid m = 0),$$

where  $\mathbb{P}(\cdot)$  is taken jointly over all randomness in Exp. 1.

Note that membership advantage is an average-case metric, as opposed to worst-case metrics considered by other privacy frameworks such as differential privacy. Our analysis is thus focused on average-case privacy leakage, and we do not provide worst-case guarantees.

## 3 Theoretical Results

### 3.1 Optimal Membership Inference Via Hypothesis Testing

In this paper, rather than studying current MI attacks, we study the optimal MI adversary: the adversary that maximizes membership advantage. As such, our analysis and results are not restricted to the current known attack strategies, which are constantly evolving, and instead serve as upper bounds for the performance of any MI attack, now or in the future. The following proposition supplies the theoretically optimal MI adversary.

**Proposition 3.1.** *The adversary that maximizes membership advantage is:*

$$A^*(\mathbf{x}_0, f_S(\mathbf{x}_0)) = \begin{cases} 1 & \text{if } P(\hat{y}_0 \mid m = 1, \mathbf{x}_0) > P(\hat{y}_0 \mid m = 0, \mathbf{x}_0), \\ 0 & \text{otherwise,} \end{cases}$$

where  $\hat{y}_0 = f_S(\mathbf{x}_0)$  and  $P$  denotes the distribution function for  $\hat{y}_0$  over the randomness in the membership inference experiment conditioned on  $\mathbf{x}_0$ .

*Proof.* Conditioned on  $m = 0$ , we have that  $(\mathbf{x}_0, y_0)$  is drawn from  $\mathcal{D}$ . Conditioned on  $m = 1$ , we have that  $(\mathbf{x}_0, y_0)$  is an element chosen randomly from  $S$ , whose elements are themselves drawn from  $\mathcal{D}$ . Thus, in both the  $m = 0$  and  $m = 1$  cases,  $\mathbf{x}_0$  has the same distribution. We thus have:

$$\begin{aligned} A^* &= \arg \max_A \text{Adv}(A) \\ &= \arg \max_A \mathbb{P}(A((\mathbf{x}_0, \hat{y}_0)) = 1 \mid m = 1) - \mathbb{P}(A((\mathbf{x}_0, \hat{y}_0)) = 1 \mid m = 0) \\ &= \arg \max_A \mathbb{E}_{\mathbf{x}_0} [\mathbb{P}(A((\mathbf{x}_0, \hat{y}_0)) = 1 \mid m = 1, \mathbf{x}_0) - \mathbb{P}(A((\mathbf{x}_0, \hat{y}_0)) = 1 \mid m = 0, \mathbf{x}_0)] \\ &= \arg \max_A \mathbb{E}_{\mathbf{x}_0} \left[ \int_{\mathbb{R}} \mathbb{1}_{A(\mathbf{x}_0, \hat{y}_0)=1} (P(\hat{y}_0 \mid m = 1, \mathbf{x}_0) - P(\hat{y}_0 \mid m = 0, \mathbf{x}_0)) dP \right], \end{aligned}$$

where in the third line, the randomness over  $\mathbf{x}_0$  is removed from the probability. To maximize the integral, we set  $A(\mathbf{x}_0, \hat{y}_0) = 1$  if  $P(\hat{y}_0 \mid m = 1, \mathbf{x}_0) - P(\hat{y}_0 \mid m = 0, \mathbf{x}_0) > 0$  and 0 otherwise.  $\square$

As observed by [10, 29], the optimal adversary performs a hypothesis test with respect to the posterior probabilities, with the two hypotheses being:

$$H_0 : \mathcal{S} \sim \mathcal{D}^n, (\mathbf{x}_0, y_0) \sim \mathcal{D} \quad \text{and} \quad H_1 : \mathcal{S} \sim \mathcal{D}^n, (\mathbf{x}_0, y_0) \sim S,$$

so that maximizing membership advantage is achieved by performing a likelihood ratio test.

### 3.2 Linear Regression with Gaussian Data

We begin our study of overparameterization's effect on MI with linear regression with Gaussian data. We find that in the sufficiently overparameterized regime, increasing the number of parameters increases vulnerability to MI. We denote by  $n, p, D$  the number of data points, number of model parameters, and data dimensionality, respectively. Let  $n, p, D \in \mathbb{Z}^+$  be given such that  $p \leq D$ , and let  $\sigma > 0$  be given. Consider a  $D$ -dimensional random vector  $\beta \sim \mathcal{N}(0, \frac{1}{D} \mathbf{I}_D)$ , representing the true coefficients. We consider data points  $(\mathbf{x}_i, y_i)$  where  $\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_D)$  and  $y_i = \mathbf{x}_i^\top \beta + \epsilon_i$ , where  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ . Sampling  $n$  data points, we denote by  $\mathbf{X}$  the  $n \times D$  matrix whose  $i^{\text{th}}$  row is  $\mathbf{x}_i^\top$  and by  $\mathbf{y}$  the  $n$ -dimensional vector of elements  $y_i$ . Further, let  $\mathbf{X}_p$  be the  $n \times p$  matrix containing the first  $p$  columns of  $\mathbf{X}$ . Least squares linear regression finds the minimum-norm  $p$ -dimensional vector  $\hat{\beta}$  that minimizes  $\|\mathbf{y} - \mathbf{X}_p \beta\|_2^2$ , which is  $\hat{\beta} = \mathbf{X}_p^\dagger \mathbf{y}$ , where  $\mathbf{X}_p^\dagger$  denotes the Moore-Penrose inverse of  $\mathbf{X}_p$ . Then, for a given feature vector  $\mathbf{x}_0$ , the model predicts  $\hat{y}_0 = \mathbf{x}_{0,p}^\top \hat{\beta}$ , where  $\mathbf{x}_{0,p}$  is the vector containing the first  $p$  elements of  $\mathbf{x}_0$ . In this setting, we can derive the membership advantage of the optimal adversary as a function of the number of parameters  $p$  used by the model.

**Theorem 3.2.** *Let  $n, p, D \in \mathbb{Z}^+$  be given such that  $n + 1 < p \leq D$ . Let  $\mathbf{x}_0$  be a given  $D$ -dimensional vector. Let  $\beta \sim \mathcal{N}(0, \frac{1}{D} \mathbf{I}_D)$ ,  $\mathbf{x}_i \sim \mathcal{N}(0, \mathbf{I}_D)$ ,  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ , and  $y_i = \mathbf{x}_i^\top \beta + \epsilon_i$  for  $i \in \{1, 2, \dots, n\}$ . Let  $m$  be a variable whose value is either 0 or 1 such that, if  $m = 1$ ,  $\mathbf{x}_k$  is set to  $\mathbf{x}_0$  for  $k$  chosen uniformly at random from  $1, 2, \dots, n$ . Let  $\mathbf{X}$  be the  $n \times D$  matrix whose rows are  $\mathbf{x}_1^\top, \mathbf{x}_2^\top, \dots, \mathbf{x}_n^\top$ . Finally, let  $\hat{y}_0 = \mathbf{x}_{0,p}^\top \mathbf{X}_p^\dagger \mathbf{y}$  where  $\mathbf{y}$  is the  $n$ -dimensional vector with elements  $y_i$ . Then, as  $n, p, D \rightarrow \infty$  such that  $\frac{p}{n} \rightarrow \gamma \in (1, \infty)$ , we have:*

$$(\hat{y}_0 \mid m = 0, \mathbf{x}_0) \sim \mathcal{N}(0, \sigma_0^2) \quad \text{and} \quad (\hat{y}_0 \mid m = 1, \mathbf{x}_0) \sim \mathcal{N}(0, \sigma_1^2) \quad (1)$$

where

$$\sigma_0^2 := \left(\frac{n}{p}\right) \left(\frac{1}{D} + \frac{1 + \sigma^2 - \frac{p}{D}}{p - n - 1}\right) \|\mathbf{x}_{0,p}\|^2 \quad \text{and} \quad \sigma_1^2 := \sigma^2 + \frac{\|\mathbf{x}_0\|^2}{D}. \quad (2)$$

Consider the case when  $\sigma_1 > \sigma_0$ . The resulting optimal membership inference algorithm  $A^*$  is

$$A^*(\mathbf{x}_0, \hat{y}_0) := \mathbb{1}[\hat{y}_0^2 > \alpha^2] \quad \text{where} \quad \alpha := \sqrt{\frac{\sigma_0^2 \sigma_1^2 \log\left(\frac{\sigma_1^2}{\sigma_0^2}\right)}{\sigma_1^2 - \sigma_0^2}}. \quad (3)$$

with membership inference advantage  $\text{Adv}$ :

$$\text{Adv}(A^*) = \mathbb{E}_{\mathbf{x}_0} \left[ 2 \left\{ \Phi\left(\frac{\alpha}{\sigma_0}\right) - \Phi\left(\frac{\alpha}{\sigma_1}\right) \right\} \right], \quad (4)$$

where  $\Phi(\cdot)$  denotes the CDF of the standard normal, and  $\alpha, \sigma_0, \sigma_1$  are conditioned on  $\mathbf{x}_0$ .

*Remark 3.3.* The case where  $\sigma_1 < \sigma_0$  occurs when  $\gamma$  is small ( $p \approx n$ ). The same membership advantage result holds in this, reversing the roles of  $\sigma_0$  and  $\sigma_1$  in  $\text{Adv}(A^*)$  in Eqs. (3) and (4) and reversing the inequality in Eq. (3).

*Remark 3.4.* The above result holds using the asymptotic distributions as  $n, p, D \rightarrow \infty$ . In Lemma C.1 in the Appendix, we derive the non-asymptotic distributions for the predictions of the minimum-norm least squares interpolator, though they cannot be written in closed form.

In Theorem 3.2, Eq. (1) shows that the posterior distributions of the outputs for test points (when  $m = 0$ ) and training points (when  $m = 1$ ) are both 0-mean Gaussians but with variances  $\sigma_0^2$  and  $\sigma_1^2$  respectively given in Eq. (2). Recall from Prop. 3.1 that the optimal MI adversary is a likelihood ratio test (LRT) between the distribution of the model’s output for new test points (when  $m = 0$ ) and that for reused training points (when  $m = 1$ ). This is reflected in Eq. (3) where we note that  $\alpha$  is the standard sufficient statistic for an LRT between two 0-mean Gaussian distributions with deviations  $\sigma_0$  and  $\sigma_1$ . Finally, we compute the membership advantage of this adversary in Eq. (4) by taking an expectation over the random draw of  $\mathbf{x}_0$ , as defined in Experiment 1, noting that  $\sigma_0, \sigma_1$ , and  $\alpha$  are all functions of  $\mathbf{x}_0$ . Membership advantage is defined to be the difference between the true and false positive probabilities in Defn. 2.1. This difference is given by  $\Phi(\alpha/\sigma_0) - \Phi(\alpha/\sigma_1)$  for each side of the Gaussian distributions, leading to the expression in Eq. (4).

To understand the implication of the result for MI, first observe that when  $m = 1$  and  $\mathbf{x}_0$  is a training point that is memorized by the model, the variance of the model’s output is equal to the variance of the measurement  $y_0$  itself independent of  $p$ . On the contrary, when  $m = 0$  and  $\mathbf{x}_0$  is a test point, the variance of the model’s output is *decreasing* with  $p$  (Figure 1a). Hence, though the output distribution means stay the same, as the variance for the  $m = 0$  case decreases far below that of the  $m = 1$  case, an LRT can easily distinguish these two distributions. In the extreme case, suppose  $p, D \gg n$ , then  $\sigma_0^2 \rightarrow 0$ , while  $\sigma_1^2$  remains a nonzero value. Hence  $\text{Adv}(A^*) \rightarrow 1$ .

We confirm our derivations of the output distributions numerically in Figure 1 where we plot empirical histograms of the predicted distributions from Eq. (1) by repeatedly computing the minimum norm least squares solution over 20,000 independent trials<sup>1</sup>. We also plot the standard deviations from Eq. (2) in Fig. 1a. Visually, from the distributions in Figure 1, given  $f(\mathbf{x}_0)$ , MI reduces to identifying whether it is more likely that a sample came from the blue distribution or from the orange with the adversary simply predicting the more likely outcome. These distributions intersect at  $\pm\alpha$ , so it is sufficient to compare  $\hat{y}_0^2$  to  $\alpha^2$  to perform the LRT.

In Figure 2a, we plot the membership advantage for  $\gamma > 1$  using Eq. (4). Since the difference between the variances of the model’s output when  $m = 0$  and  $m = 1$  increases with  $\gamma$  for  $\gamma > 2$  (cf., Fig. 1a), it becomes easier to distinguish the two distributions. This results in an increase in the membership advantage. The

<sup>1</sup>We detail additional experimental details and the computational hardware in Appendix D.

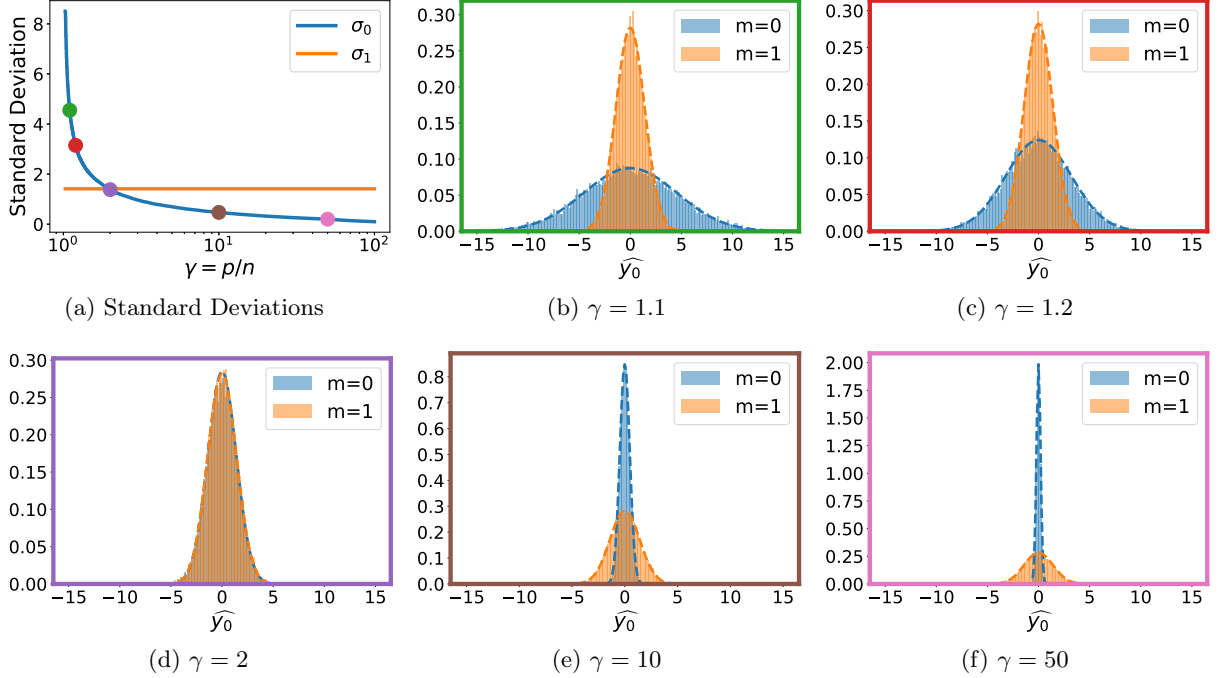


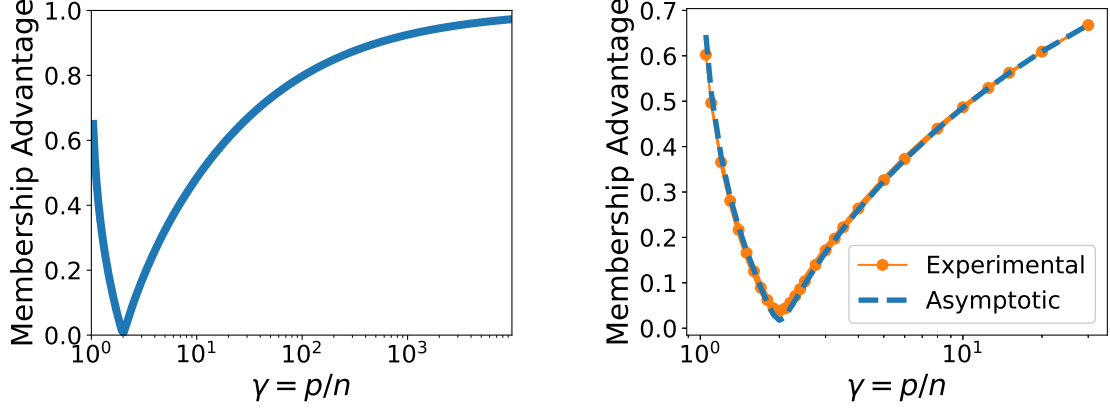
Figure 1: In the highly overparametrized regime, increasing the number of parameters  $p$  yields more distinguishable posterior distributions. (a) Standard deviations of model outputs  $\hat{y}_0$  for minimum-norm least squares as a function of parameterization for  $\gamma > 1$ . (b-f) The Gaussian distributions in Eq. (1) (broken lines), as well as empirical histograms of 20,000  $\hat{y}_0$  samples for different  $\gamma$  values with  $n = 400$ ,  $D = 20,000$ ,  $\sigma = 1$ . The prediction variance for  $m = 1$  stays constant, while the variance for  $m = 0$  decreases with increased parameterization, making the distributions easier to distinguish.

initial decrease in membership advantage when  $\gamma \leq 2$  is a consequence of  $\sigma_1 \leq \sigma_0$  in this regime as shown in Fig. 1a. Since  $\sigma_0$  is decreasing in  $p$ , initially, this decrease makes the train and test output distributions harder to distinguish, leading to lower membership advantage and 0 advantage for  $\gamma = 2$ . However, for  $\gamma > 2$ ,  $\sigma_0$  decreases past  $\sigma_1$  and membership advantage approaches 1 as  $\gamma \rightarrow \infty$ . In practice, the  $\gamma < 2$  regime (only slightly overparametrized) is less interesting since models suffer larger generalization error in this setting (cf., Theorem 1 of [23]) and are rarely (if ever) used in practice. A key takeaway is that for linear regression in the Gaussian data setting, extreme overparameterization increases the vulnerability of a machine learning model to MI.

While Theorem 3.2 operates in the asymptotic regime, we empirically approximate membership advantage for  $n = 100$  and  $D = 3,000$  by approximating the posterior distributions with histograms of samples. As we see in Figure 2b, the asymptotic formula agrees very closely with the non-asymptotic experiment. Experimental details are provided in Appendix D.

## 4 Mitigating Membership Inference Attacks

Next, we extend our analysis from the previous section towards two methods commonly used for preserving privacy: regularization and noise addition. We present two key results. First, we show that for the overparametrized Gaussian data setting, ridge regularization actually *increases* membership advantage and is thus detrimental to privacy. Second, we show that the privacy-utility trade-off induced by reducing the number of parameters of a linear regression model with Gaussian data is equivalent to that of adding independent Gaussian noise to the output of a model that uses all available features.



(a) Asymptotic membership advantage

(b) Empirical verification for non-asymptotic case

Figure 2: Increasing the number of parameters increases membership advantage. (a) Theoretical membership advantage for linear regression with Gaussian data (Eq. 4) as a function of the number of parameters for 100 sampled  $\mathbf{x}_0$ 's,  $n = 10^3$ ,  $D = 10^7$ . (b) We empirically approximate the membership advantage averaged over 20 sampled  $\mathbf{x}_0$ 's for  $n = 100$ ,  $D = 3,000$ , and  $\sigma = 1$  by estimating the two posterior distributions for each  $\mathbf{x}_0$  using empirical histograms with 100,000 samples. We plot it alongside the theoretical asymptotic membership advantage, showing a close agreement between the two.

#### 4.1 Ridge-Regularized Linear Regression

We analyze membership inference in *ridge-regularized* linear regression in the same setting as Section 3.2 except that, now, the estimate is  $\hat{\beta}_\lambda = (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^\dagger \mathbf{X}_p^\top y$ , where  $\lambda$  is a regularization parameter. Larger values of  $\lambda$  yield greater regularization, and  $\lambda \rightarrow 0$  reduces to the case of Section 3.2. Regularization is a common method to reduce overfitting and has thus been proposed in previous works as a defense mechanism to protect models from MI attacks [13, 30].

A surprising observation resulting from our analysis is that, in the highly overparameterized regime ( $\gamma \gg 1$ ), ridge regularization actually *increases* the model's vulnerability to MI attacks for linear regression with Gaussian data. In Theorem 4.1, we derive an analogous result for the distributions of the ridge-regularized predictions as Theorem 3.2 demonstrated for the unregularized case.

**Theorem 4.1. Membership advantage for ridge-regularized linear regression.** *Consider the same setup as in Theorem 3.2, but now let  $\hat{\beta}_\lambda = (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top$  for some  $\lambda > 0$ . Then, as  $n, p, D \rightarrow \infty$  such that  $\frac{p}{n} \rightarrow \gamma \in (1, \infty)$ , we have:*

$$(\hat{y}_0 \mid m = 0, \mathbf{x}_0) \sim \mathcal{N}(0, \sigma_{0,\lambda}^2) \quad \text{and} \quad (\hat{y}_0 \mid m = 1, \mathbf{x}_0) \sim \mathcal{N}(0, \sigma_{1,\lambda}^2),$$

where

$$\begin{aligned}
\sigma_{0,\lambda}^2 &:= \frac{g'(-\lambda)\gamma}{(1+g(-\lambda)\gamma)^2} \left( \sigma^2 + 1 - \frac{p}{D} \right) \frac{\|\mathbf{x}_{0,p}\|_2^2}{p} + (1 - 2\lambda g(-\lambda) + \lambda^2 g'(-\lambda)) \frac{\|\mathbf{x}_{0,p}\|_2^2}{D} \\
\sigma_{1,\lambda}^2 &:= \left[ \left( \frac{\lambda^2}{(\lambda + \gamma g(-\lambda))(\lambda + \gamma \frac{\|\mathbf{x}_{0,p}\|_2^2}{p} g(-\lambda))} \right)^2 \gamma g'(-\lambda) \frac{\|\mathbf{x}_{0,p}\|_2^2}{p} \right] \left( \sigma^2 + 1 - \frac{p}{D} \right) \\
&\quad + \left( \frac{\gamma g(-\lambda) \frac{\|\mathbf{x}_{0,p}\|_2^2}{p}}{1 + \gamma g(-\lambda) \frac{\|\mathbf{x}_{0,p}\|_2^2}{p}} \right)^2 \left( \sigma^2 + \frac{\|\mathbf{x}_{0,p}\|_2^2}{D} \right) \\
&\quad + \left( 1 - \frac{2\lambda g(-\lambda)}{1 + \gamma \frac{\|\mathbf{x}_{0,p}\|_2^2}{p} g(-\lambda)} + \frac{\lambda^2 g'(-\lambda)}{\left( 1 + \gamma \frac{\|\mathbf{x}_{0,p}\|_2^2}{p} g(-\lambda) \right)^2} \right) \frac{\|\mathbf{x}_{0,p}\|_2^2}{D}, \\
g(-\lambda) &:= \frac{-(1 - \gamma + \lambda) + \sqrt{(1 - \gamma + \lambda)^2 + 4\gamma\lambda}}{2\gamma\lambda}.
\end{aligned}$$

Furthermore, in the case when  $\sigma_{1,\lambda} > \sigma_{0,\lambda}$  and defining:

$$\alpha_\lambda = \sqrt{\frac{\sigma_{0,\lambda}^2 \sigma_{1,\lambda}^2 \log\left(\frac{\sigma_{1,\lambda}^2}{\sigma_{0,\lambda}^2}\right)}{\sigma_{1,\lambda}^2 - \sigma_{0,\lambda}^2}},$$

the optimal membership inference advantage is then:

$$Adv(A_\lambda^*) = \mathbb{E}_{\mathbf{x}_0} \left[ 2 \left\{ \Phi\left(\frac{\alpha_\lambda}{\sigma_{0,\lambda}}\right) - \Phi\left(\frac{\alpha_\lambda}{\sigma_{1,\lambda}}\right) \right\} \right].$$

*Remark 4.2.* The above result holds using the asymptotic distributions as  $n, p, D \rightarrow \infty$ . In Lemma C.3, we derive the non-asymptotic distributions for the predictions of the ridge-regularized least squares estimator, though they cannot be written in closed form.

The membership advantage  $Adv(A_\lambda^*)$  achieved by the optimal adversary  $A_\lambda^*$  of the  $\lambda$ -regularized model is an increasing function in  $\lambda$  when  $\gamma \gg 1$ . To visualize this, we plot the theoretical membership advantages for 100 sampled  $\mathbf{x}_0$ 's (each with iid standard normal elements) for differing regularization strengths in Figure 3 with the same setting as in Figure 2a. In Figure 3a, we plot membership advantage as a function of the overparameterization ratio  $\gamma$  for a few  $n\lambda$  values. Conversely, In Figure 3b we plot membership advantage versus regularization  $n\lambda$  for a few different values of  $\gamma$ . In particular, note how in both subplots, increasing  $\lambda$  never decreases the membership advantage. This observation has been made empirically (but not analytically) for techniques that have similar regularizing effects, such as dropout [15], ensembling [22], and weight decay [30].

Intuitively, the reason ridge regularization increases MI vulnerability is because while it decreases the variance of the model's output on training points, it also significantly decreases the variance for test points such that the two output distributions become more distinguishable. In Figure 6 in the appendix, we plot the gap between the variances for the  $m = 0$  and  $m = 1$  cases for different regularization amounts  $\lambda$  to visualize this effect.

## 4.2 Noise Addition vs. Feature Reduction

A consequence of Thm. 3.2 is that, in the Gaussian data setting, one can reduce vulnerability to MI attacks by simply *decreasing* the number of parameters/features. However, this comes at the cost of decreased generalization performance (utility) due to the ‘‘double descent’’ effect [6, 7, 23], wherein generalization error



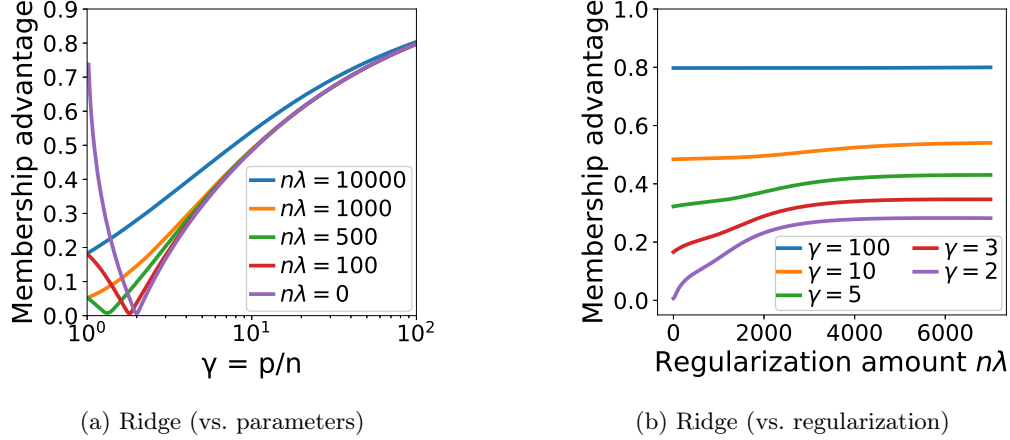


Figure 3: For highly overparameterized models, ridge regularization does not decrease membership advantage (MA) in the Gaussian data setting. (a) Theoretical MA vs. overparameterization  $\gamma$  for a few ridge regularization strengths  $\lambda$ . Stronger ridge regularization harms privacy (increased MA). (b) Theoretical MA vs.  $\lambda$  for a few  $\gamma$ 's, explicitly showing how increasing  $\lambda$  increases MA. For this analysis, we set  $n = 10^3$ ,  $D = 10^7$ , and  $\sigma = 1$ .

decreases with increased overparameterization. An alternative and popular method to increase a model's privacy is adding independent noise to the model output [24, 31], but this also decreases generalization performance. Interestingly, we show in this subsection that for Gaussian data, the privacy-utility trade-offs induced by both feature reduction and noise addition are actually equivalent.

**Feature reduction:** To characterize the privacy-utility trade-off of feature reduction, we first need the generalization error for a given number of parameters, provided in Corollary 2.2 of [23]:

**Proposition 4.3.** (Adapted from Corollary 2.2 of [23]). *In the same setting as Theorem 3.2, for  $\mathbf{x}_0 \sim \mathcal{N}(0, \mathbf{I}_D)$ , the generalization error is given by*

$$\mathbb{E}[(y - \hat{\beta}^\top \mathbf{x}_{0,p})^2] = 1 + \sigma^2 + n \left( \frac{1 + \sigma^2 - \frac{p}{D}}{p - n - 1} - \frac{1}{D} \right).$$

Alternatively, this expression can be written as  $1 + \sigma^2 + \mathbb{E}_{\mathbf{x}_0}[\sigma_0^2] - 2\frac{n}{D}$ , with  $\sigma_0^2$  as in Theorem 3.2.

**Noise addition:** Next, we consider noise addition performed by perturbing the  $m = 0$  model output with independent noise before releasing it:  $\hat{y}_0 = \mathbf{x}_{0,p}^\top \hat{\beta} + \bar{\epsilon}$  where  $\bar{\epsilon} \sim \mathcal{N}(0, \bar{\sigma}^2)$ . Recall that MI is possible for overparameterized linear regression models because the variance of the output of test points ( $m = 0$ ) is lower than the variance of the output for training points ( $m = 1$ ) (Fig. 1). Thus, adding random noise to the model output when a test point  $\mathbf{x}_0$  is queried can make the output distributions harder to distinguish. In the following lemma, we compute the membership advantage and generalization error of noise addition using the full set of features  $p = D$ .

**Lemma 4.4.** *Consider the setup of Theorem 3.2 restricted to  $p = D$ ,  $\sigma_1 > \sigma_0$ , and  $\gamma > 1$ , and suppose independent 0-mean, variance  $\bar{\sigma}^2$  Gaussian noise is added to the outputs when  $m = 0$  ( $\mathbf{x}_0$  is not a member of the training set). Then,  $(\hat{y}_0 \mid m = 0, \mathbf{x}_0) \sim \mathcal{N}(0, \sigma_0^2 + \bar{\sigma}^2)$ . The optimal adversary is*

$$A_{\bar{\sigma}}^*(\mathbf{x}_0, \hat{y}_0) := \mathbb{1}[\hat{y}_0^2 > \alpha_{\bar{\sigma}}^2] \quad \text{where} \quad \alpha_{\bar{\sigma}} := \sqrt{\frac{(\sigma_0^2 + \bar{\sigma}^2)\sigma_1^2 \log\left(\frac{\sigma_1^2}{(\sigma_0^2 + \bar{\sigma}^2)}\right)}{\sigma_1^2 - (\sigma_0^2 + \bar{\sigma}^2)}}.$$

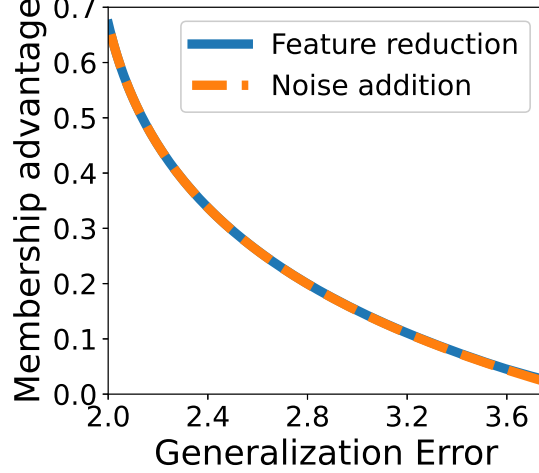


Figure 4: We plot the privacy-utility trade-off obtained when tuning the number of parameters (feature reduction; blue) and when adding independent noise to the output of a model that uses all available parameters (noise addition; orange) and demonstrate that the two trade-offs are essentially equivalent. We use  $n = 100$ ,  $D = 3000$ , and  $\sigma = 1$  for this analysis.

*Its membership inference advantage is*

$$\text{Adv}(A_{\bar{\sigma}}^*) = \mathbb{E}_{\mathbf{x}_0} \left[ 2 \left\{ \Phi \left( \frac{\alpha_{\bar{\sigma}}}{\sqrt{\sigma_0^2 + \bar{\sigma}^2}} \right) - \Phi \left( \frac{\alpha_{\bar{\sigma}}}{\sigma_1} \right) \right\} \right]. \quad (5)$$

*Additionally, for  $\mathbf{x}_0 \sim \mathcal{N}(0, \mathbf{I}_D)$ , the generalization error incurred is*

$$\mathbb{E}[(y - \hat{\beta}^\top \mathbf{x}_0)^2] = 1 + \sigma^2 + n \left( \frac{\sigma^2}{D - n - 1} - \frac{1}{D} \right) + \bar{\sigma}^2. \quad (6)$$

The terms in Lemma 4.4 are similar to those in Theorem 3.2 except for the added dependence on the added noise’s variance  $\bar{\sigma}^2$  in the case that  $m = 0$  ( $\mathbf{x}_0$  is not a member). Increasing noise variance  $\bar{\sigma}^2$  decreases membership advantage (Eq. (5)) at the cost of increased generalization error (Eq. (6)). Indeed, it is possible to add sufficient noise such that  $\sigma_1^2 = \sigma_0^2 + \bar{\sigma}^2$ , rendering the membership advantage 0, though possibly at the cost of impermissible generalization performance.

In Figure 4, we plot the membership advantage vs. generalization error trade-offs for both feature reduction (blue) and noise addition (orange). The plots follow the setting of Fig. 2b. For the blue curve, we employ the expressions in Theorem 3.2 and Proposition 4.3 while varying the overparameterization ratio  $\gamma$ . For the orange curve, we use Lemma 4.4, using all available features ( $p = D$ ) while varying the noise variance  $\bar{\sigma}^2$ . Fig. 4 shows that the trade-offs induced by both feature reduction and noise addition are essentially equivalent.

## 5 More Complex Models

In this section, we present three more complex data models wherein we empirically observe increased overparameterization leading to increased MI vulnerability.

For each data model, we perform the following experiment. We first sample a  $\mathbf{x}_0$  vector, which is the data point we wish to perform MI on. Then we sample a training dataset  $\mathbf{X}$ , measurements  $\mathbf{y}$ , as well as other random elements according to the data model. To obtain an  $m = 0$  prediction, we learn a model on  $\mathbf{X}$  that we

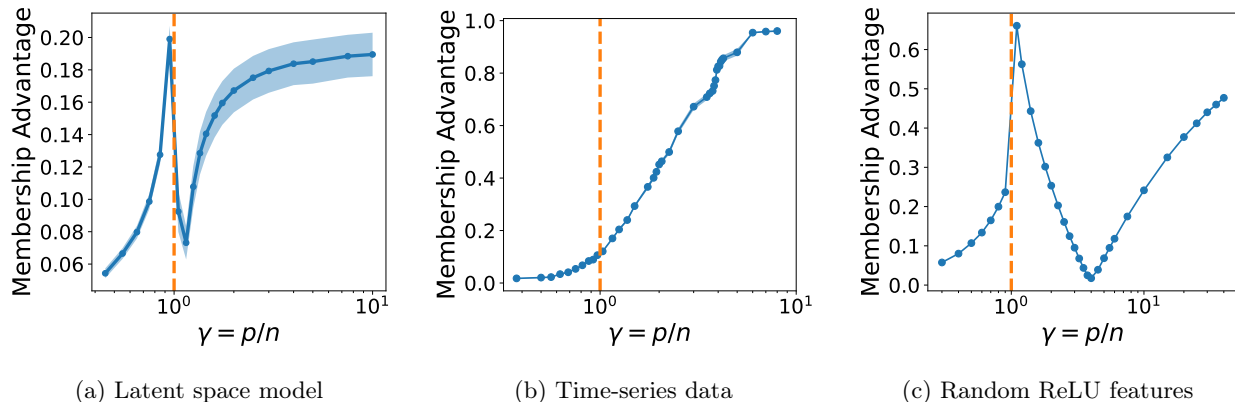


Figure 5: Empirically measured membership advantage vs. parameterization for various data models detailed in Section 5. The dashed line at  $\gamma = 1$  divides the underparameterized and overparameterized regions. For all settings, when sufficiently overparameterized, increasing the number of parameters increases vulnerability to membership inference attacks. (a) Latent space model with  $n = 200$  where  $p$  covariates of  $d = 20$  latent features are observed. (b) Regression over  $n = 128$  time samples of a linear combination of  $D = 1024$  Fourier features. (c) Nonlinear random ReLU features with  $n = 100$ ,  $D = 5000$ , and  $\sigma = 1$ .

then apply on  $\mathbf{x}_0$ . To obtain an  $m = 1$  prediction, we first replace a row of  $\mathbf{X}$  with  $\mathbf{x}_0$  and the corresponding element of  $\mathbf{y}$  with  $y_0$  before learning the model and applying it on  $\mathbf{x}_0$ . Keeping  $\mathbf{x}_0$  fixed throughout the experiment and resampling all other random data (such as  $\mathbf{X}$ ) many times, we collect a large set of  $m = 0$  and  $m = 1$  prediction samples. We build a histogram of these samples by assigning them into fine discrete bins to obtain approximations of the conditional densities  $\hat{P}(\hat{y}_0 | m = 1, \mathbf{x}_0)$  and  $\hat{P}(\hat{y}_0 | m = 0, \mathbf{x}_0)$  needed for the optimal adversary (cf., Prop 3.1). To approximate membership advantage, we sum up the differences  $\hat{P}(\hat{y}_0 | m_1, \mathbf{x}_0) - \hat{P}(\hat{y}_0 | m_0, \mathbf{x}_0)$  over all the histogram bins where  $\hat{P}(\hat{y}_0 | m_1, \mathbf{x}_0) > \hat{P}(\hat{y}_0 | m_0, \mathbf{x}_0)$ . We repeat this experiment 20 times, each with a newly sampled  $\mathbf{x}_0$ , and plot the means (as points) and the estimated standard errors (as shaded error regions) of the membership advantage values across the 20 experiments in Figure 5. We next discuss the data models in detail.

**Linear Regression on Latent Space Model:** We first consider the latent space model from [3], where the output variable  $y_i$  is a noisy linear function of a data point’s  $d$  latent features  $\mathbf{z}_i$ , but one only observes a vector  $\mathbf{x}_i$  containing  $p \geq d$  covariates rather than the direct features  $\mathbf{z}_i$ . Let  $\mathbf{Z}$  be an  $n \times d$  matrix where each row is the vector of latent features for an observation. We then have:

$$\mathbf{y} = \mathbf{Z}\beta + \epsilon, \quad \mathbf{X}_{i,j} = \mathbf{w}_j^\top \mathbf{z}_i + u_{i,j}, \quad \hat{y}_0 = \mathbf{x}_0^\top \mathbf{X}_p^\dagger \mathbf{y},$$

where  $\mathbf{w}_j$  is a  $d$ -dimensional vector, and  $u_{i,j}$  is a noise term. In this experiment, we set  $n = 200$ ,  $d = 20$ , and vary  $p$ . For each experiment, we sample a single  $\mathbf{x}_0 \sim \mathcal{N}(0, \mathbf{I}_d)$  and a single set of  $\mathbf{w}_j$  vectors, each from  $\mathcal{N}(0, \mathbf{I}_d)$ , and keep them fixed. We leave the other variables random with the following distributions:  $\mathbf{z}_i \sim \mathcal{N}(0, \mathbf{I}_d)$ ,  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ ,  $\beta \sim \mathcal{N}(0, \frac{1}{d}\mathbf{I}_d)$  and  $u_{i,j} \sim \mathcal{N}(0, 1)$ . In Fig. 5a, we plot the empirical membership advantage values, which increase with the number  $p$  of features in the overparameterized regime.

**Noise-Free Time-Series Regression:** In this experiment, we consider a model that aims to interpolate a time-series signal using frequency components as the features. For example, consider a patient who visits a hospital at irregular times  $t_i$  to get their blood glucose level measured. After obtaining a number of measurements taken over time, the hospital fits a time-series signal representing the patient’s blood glucose level at any time. Using this learned model, an adversary wishes to identify whether the patient visited the hospital at a particular time  $t_0$ , that is, if  $t_0$  is one such time point included in the hospital’s training dataset for learning the patient’s blood glucose level.

To formalize this, we fix  $D = 1024$  as the number of frequency components each signal contains and let  $M = 2D - 1$ . Let  $\mathbf{W}$  be the  $D \times D$  matrix whose elements  $\mathbf{W}_{kl} = \frac{1}{\sqrt{M+2}} \cos\left(\frac{2\pi kl}{M}\right)$ . That is, each

column of the matrix is half a period of an  $M$ -dimensional discrete cosine with frequency  $\frac{2\pi k}{M}$ . Sample a  $\beta \sim \mathcal{N}(0, \frac{1}{D}\mathbf{I}_D)$ , and let  $\mathbf{z} = \mathbf{W}\beta$  denote the true length- $D$  signal. Thus, the signal is a random linear combination of cosines. We randomly select  $n = 128$  indices from  $1, 2, \dots, D$  uniformly without replacement, and let  $\mathbf{X}$  be the  $n \times D$  matrix whose rows are the rows of  $\mathbf{W}$  at these selected indices. Then,  $\mathbf{y} = \mathbf{X}\beta$  is the signal observed at the randomly selected  $n$  indices. The regressor learns  $\hat{\beta} = \mathbf{X}_p^\dagger \mathbf{y}$  and then predicts the signal at any other time point  $t_0$  as  $\hat{y}_0 = \mathbf{x}_{0,p}^\top \hat{\beta}$ , where  $\mathbf{x}_0$  is row  $t_0$  of  $\mathbf{W}$ . Thus, identifying whether  $t_0$  was a time point in the training dataset is equivalent to identifying if  $\mathbf{x}_0$  was in  $\mathbf{X}$ . The membership advantage values for this task, plotted in Figure 5b, increases with the number  $p$  of frequency components included in the model.

**Random ReLU Features:** We next consider a nonlinear data model based on Random ReLU feature networks [32, 33]. Let  $\mathbf{Z}$  be a random  $n \times D$  matrix whose elements are iid standard normal. Let  $\mathbf{V}$  be a random  $D \times p$  matrix whose rows are sampled iid from the surface of the unit sphere in  $\mathbb{R}^p$ . Let  $\mathbf{X} = \max(\mathbf{Z}\mathbf{V}, 0)$ , where the max is taken elementwise. The target variables are given by  $\mathbf{y} = \mathbf{Z}\beta + \sigma\epsilon$ , where  $\beta \sim \mathcal{N}(0, \frac{1}{D}\mathbf{I}_D)$  and  $\epsilon \sim \mathcal{N}(0, \mathbf{I}_n)$ . Finally, for the data point  $\mathbf{x}_0$ , let its prediction be  $\hat{y}_0 = \mathbf{x}_0^\top \mathbf{X}^\dagger \mathbf{y}$ . We plot the membership advantage in Figure 5c with  $n = 100$ ,  $D = 5000$ , and  $\sigma = 1$ . We again observe that in the highly overparametrized regime, membership advantage increases with parameters.

## 6 Discussion and Conclusions

We have shown theoretically for (regularized) linear regression with Gaussian data and empirically for more complex models (latent space regression, time-series regression using Fourier components, and random ReLU features) that increasing the number of model parameters renders them more vulnerable to membership inference attacks. Thus, while overparameterization may be attractive for its robust generalization performance, one must proceed with caution to ensure the learned model does not lead to unintended privacy leakages.

More speculatively, we hypothesize that the same overparameterization/vulnerability tradeoff should exist in many machine learning models (e.g., deep networks) beyond those we have studied. Intuitively, the output of a model that achieves zero training error but generalizes well must i) fit to any noise (e.g. additive Gaussian noise) present in the training data to get a perfect fit to the noisy training data but also ii) eliminate the effect of noise in the training data in predicting for unseen data to achieve good generalization. This disparate behavior towards training and non-training data points leads to different output distributions when the input is or is not among the training data and is universal for overparameterized models. Ultimately, this causes a difference in the distributions of training and test predictions that can be leveraged to perform a membership attack.

There are still many open questions in this line of research. While we have shown multiple settings where reducing the number of parameters can increase privacy, it remains to be verified that the phenomenon holds widely for other types of machine learning settings such as language tasks or large-scale image recognition. Another interesting next step would be to investigate how increased overparameterization affects privacy for models trained with privacy-preserving techniques or membership inference defense schemes other than ridge regularization. We believe the findings of our work can provide insights towards developing the next generation of privacy-preserving techniques. It is our hope that the observations and analyses in this paper take a step towards keeping sensitive training data safer in a world increasingly intertwined with machine learning.

## Acknowledgements

This work was supported by NSF grants CCF-1911094, IIS-1838177, and IIS-1730574; ONR grants N00014-18-12571, N00014-20-1-2534, and MURI N00014-20-1-2787; AFOSR grant FA9550-22-1-0060; and a Vannevar Bush Faculty Fellowship, ONR grant N00014-18-1-2047.

## A Additional Discussions

### A.1 Limitations

In this work, we focus on the optimal membership inference adversary. We study this because of how it serves as an upper bound for all other attacks and because of how it yields interpretable and fundamental theoretical results. The optimal membership inference adversary has full knowledge of the learning model's output distributions when the data point of interest is a member or non-member of the training dataset. In practice, the adversary rarely has such full knowledge, and the learning model's output distributions have to be approximated using shadow models [13], or the entire attack has to be simplified, such as with a loss threshold [9, 14]. Our study does not analyze how our results are affected by the non-optimality of these more practical attacks.

### A.2 Ethical Considerations

It is the hope of the authors that by more clearly exposing the link between membership inference vulnerability and generalization performance, researchers can make informed decisions about how to achieve the best trade-off they can for their application. That said, by studying the performance of optimal membership inference attacks, it is possible that this work will call attention to vulnerabilities in existing model architectures which may be exploited. Furthermore, in settings where privacy is absolutely crucial, such as in medical applications, additional care should be taken to guard privacy beyond the guarantees of this work.

## B Proofs

### B.1 Tools for Asymptotic Analysis

The following lemmas are used in the proofs of Theorems 3.2 and 4.1. We begin with the following lemma, which is a generalized version of the Marchenko-Pastur theorem [34–36].

**Lemma B.1.** *Let  $\mathbf{X}_n \in \mathbb{R}^{n \times p}$  be a sequence of random matrices with i.i.d.  $\mathcal{N}(0, 1)$  entries. Consider the the sample covariance matrix  $\hat{\Sigma} = (1/n)\mathbf{X}_n^\top \mathbf{X}_n$ . Let  $\mathbf{C}_n \in \mathbb{R}^{p \times p}$  be a sequence of matrices such that  $\text{Tr}(\mathbf{C}_n)$  is uniformly bounded with probability one. As  $n, p \rightarrow \infty$  with  $p/n = \gamma \in (0, \infty)$ , it holds that almost surely,*

$$\text{Tr}\left(\mathbf{C}_n \left((\Sigma + \lambda \mathbf{I}_p)^{-1} - g(-\lambda) \mathbf{I}_p\right)\right) \rightarrow 0, \quad \text{Tr}\left(\mathbf{C}_n \left((\Sigma + \lambda \mathbf{I}_p)^{-2} - g'(-\lambda) \mathbf{I}_p\right)\right) \rightarrow 0$$

where  $g(\lambda)$  is the Stieltjes transform of the Marchenko-Pastur law with parameter  $\gamma$ .

We use the following Lemma in computing the asymptotic distribution of the output.

**Lemma B.2.** *Let  $\mathbf{y}_n \in \mathbb{R}^n$  be a sequence of i.i.d.  $\mathcal{N}(0, \mathbf{I}_n)$  random vectors. Also, let  $\mathbf{x}_n \in \mathbb{R}^n$  be a sequence of random vectors with spherically symmetric distribution such that  $\|\mathbf{x}_n\|_2 \xrightarrow{a.s.} \sigma$ . Further, assume that  $\mathbf{x}_n, \mathbf{y}_n$  are independent. Then  $\mathbf{x}_n^\top \mathbf{y}$  converges weakly to a zero mean gaussian with variance  $\sigma^2$ .*

*Proof.* We can write

$$\mathbf{x}_n^\top \mathbf{y}_n = \|\mathbf{x}_n\|_2 \left( \frac{\mathbf{x}_n}{\|\mathbf{x}_n\|} \right)^\top \mathbf{y}_n = \|\mathbf{x}_n\|_2 \mathbf{u}_n^\top \mathbf{y}_n$$

where  $\mathbf{u}_n \in S^{n-1}$  is uniformly distributed over the unit sphere and is independent from  $\mathbf{y}_n$ . Therefore, we can fix  $\mathbf{u}_n$  to be the first standard unit vector and the distribution of  $\mathbf{x}_n^\top \mathbf{y}_n$  is the same as  $\|\mathbf{x}_n\|_2 y_{n,1}$  where  $y_{n,1} \sim \mathcal{N}(0, 1)$ . Hence, using  $\|\mathbf{x}_n\|_2 \xrightarrow{a.s.} \sigma$ , we deduce the result.  $\square$

## B.2 Proof of Theorem 3.2

*Proof of Theorem 3.2.* Let  $\mathbf{X}_{\bar{p}}$  denote the matrix formed by removing the first  $p$  columns from  $\mathbf{X}$ , and let  $\beta_{\bar{p}}$  denote the vector formed by removing the first  $p$  elements from  $\beta$ . Recall that

$$\begin{aligned}\widehat{y}_0 \mid m = 0 &= \mathbf{x}_0^\top \mathbf{X}_p^\dagger (\mathbf{X}\beta + \epsilon) \\ &= \mathbf{x}_0^\top \mathbf{X}_p^\dagger (\mathbf{X}_p \beta_p + \eta)\end{aligned}$$

where  $\eta = \mathbf{X}_{\bar{p}}\beta_{\bar{p}} + \epsilon \sim \mathcal{N}(0, (1 + \sigma^2 - \frac{p}{D}) \mathbf{I}_n)$ . First note that the distributions of  $\mathbf{X}_p^\top \mathbf{X}_p^\dagger \mathbf{x}_0$  are  $\mathbf{X}_p^\dagger \mathbf{x}_0$  are spherically symmetric and letting  $\widehat{\Sigma} \triangleq (1/n) \mathbf{X}_p^\top \mathbf{X}_p$  and  $\mathbf{P}$  to be orthogonal projection onto row space of  $\mathbf{X}_p$  we have

$$\begin{aligned}\frac{1}{D} \|\mathbf{X}_p^\top \mathbf{X}_p^\dagger \mathbf{x}_0\|_2^2 &= \frac{1}{D} \|\mathbf{P} \mathbf{x}_0\|_2^2 = \frac{1}{D} \lim_{\lambda \rightarrow 0} \mathbf{x}_0^\top (\widehat{\Sigma} + \lambda \mathbf{I}_p)^{-1} \widehat{\Sigma} \mathbf{x}_0 \\ &= \frac{1}{D} \|\mathbf{x}_0\|_2^2 - \frac{1}{D} \lim_{\lambda \rightarrow 0} \lambda \mathbf{x}_0^\top (\widehat{\Sigma} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}_0,\end{aligned}$$

and,

$$\begin{aligned}\|\mathbf{X}_p^\dagger \mathbf{x}_0\|_2^2 &= \frac{1}{n} \lim_{\lambda \rightarrow 0} \mathbf{x}_0^\top (\widehat{\Sigma} + \lambda \mathbf{I}_p)^{-1} \widehat{\Sigma} (\widehat{\Sigma} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}_0 \\ &= \frac{1}{n} \lim_{\lambda \rightarrow 0} \mathbf{x}_0^\top \left[ (\widehat{\Sigma} + \lambda \mathbf{I}_p)^{-1} - \lambda (\widehat{\Sigma} + \lambda \mathbf{I}_p)^{-2} \right] \mathbf{x}_0.\end{aligned}$$

Thus, using Lemma B.2, both  $\frac{1}{D} \|\mathbf{X}_p^\top \mathbf{X}_p^\dagger \mathbf{x}_0\|_2^2$  and  $\|\mathbf{X}_p^\dagger \mathbf{x}_0\|_2^2$  converge to a fixed limit as  $n \rightarrow \infty$ , almost surely. Therefore, using Lemma B.2,  $\widehat{y}_0$  converges weakly to a gaussian. Now, we compute its variance. Since  $\eta$  and  $\beta$  are both zero-mean independent Gaussians and are thus orthogonal in expectation, we have by the Pythagorean theorem:

$$\mathbb{E}[\widehat{y}_0^2 \mid m = 0] = \mathbb{E}[(\mathbf{x}_0^\top \mathbf{X}_p^\dagger \mathbf{X}_p \beta_p)^2] + \mathbb{E}[(\mathbf{x}_0^\top \mathbf{X}_p^\dagger \eta)^2].$$

We start with the first term.

Note that since,  $p > n$ ,  $\mathbf{X}_p$  does not have linearly independent columns. Let  $\mathbf{P} = \mathbf{X}_p^\dagger \mathbf{X}_p$ . We have:

$$\begin{aligned}\mathbb{E}[(\mathbf{x}_0^\top \mathbf{X}_p^\dagger \mathbf{X}_p \beta_p)^2] &= \mathbb{E}[\text{Tr}(\mathbf{x}_0^\top \mathbf{P} \beta_p \beta_p^\top \mathbf{P}^\top \mathbf{x}_0)] \\ &= \mathbb{E}[\text{Tr}(\beta_p \beta_p^\top \mathbf{P}^\top \mathbf{x}_0 \mathbf{x}_0^\top \mathbf{P})] \\ &= \text{Tr}(\mathbb{E}[\beta_p \beta_p^\top \mathbf{P}^\top \mathbf{x}_0 \mathbf{x}_0^\top \mathbf{P}]) \\ &= \text{Tr}(\mathbb{E}[\beta_p \beta_p^\top] \mathbb{E}[\mathbf{P}^\top \mathbf{x}_0 \mathbf{x}_0^\top \mathbf{P}]) \\ &= \frac{1}{D} \text{Tr}(\mathbb{E}[\mathbf{P}^\top \mathbf{x}_0 \mathbf{x}_0^\top \mathbf{P}]) \\ &= \frac{1}{D} \mathbb{E}[\|\mathbf{P}^\top \mathbf{x}_0\|^2] \\ &= \frac{1}{D} \frac{n}{p} \|\mathbf{x}_0\|^2.\end{aligned}$$

In the last line, we use the same argument as in Section 2.2 of [23], using the facts that  $\mathbf{P}$  is the orthogonal projection to the row space of  $\mathbf{X}_p$  and that the Gaussian distribution is invariant to rotations.

We now consider the second term:

$$\begin{aligned}
\mathbb{E} \left[ (\mathbf{x}_0^\top \mathbf{X}_p^\dagger \eta)^2 \right] &= \left( 1 + \sigma^2 - \frac{p}{D} \right) \mathbf{x}_0^\top \mathbb{E} \left[ \mathbf{X}_p^\dagger \mathbf{X}_p^{\dagger\top} \right] \mathbf{x}_0 \\
&= \left( 1 + \sigma^2 - \frac{p}{D} \right) \mathbf{x}_0^\top \mathbb{E} \left[ \left( \mathbf{X}_p^\top \mathbf{X}_p \right)^\dagger \mathbf{X}_p^\top \mathbf{X}_p \left( \mathbf{X}_p^\top \mathbf{X}_p \right)^{\dagger\top} \right] \mathbf{x}_0 \\
&= \left( 1 + \sigma^2 - \frac{p}{D} \right) \mathbf{x}_0^\top \mathbb{E} \left[ \left( \mathbf{X}_p^\top \mathbf{X}_p \right)^\dagger \right] \mathbf{x}_0.
\end{aligned}$$

where  $\left( \mathbf{X}_p^\top \mathbf{X}_p \right)^\dagger$  has the generalized inverse Wishart distribution with expectation equal to  $\mathbb{E} \left[ \left( \mathbf{X}_p^\top \mathbf{X}_p \right)^\dagger \right] = \frac{n}{p} \frac{1}{p-n-1} \mathbf{I}_p$  (Theorem 2.1 in [37]). Thus, we have:

$$\mathbb{E} \left[ (\mathbf{x}_0^\top \mathbf{X}_p^\dagger \eta)^2 \right] = \left( \frac{n}{p} \right) \left( \frac{1 + \sigma^2 - \frac{p}{D}}{p - n - 1} \right) \|\mathbf{x}_0\|^2$$

Adding this with the result for the first term gives the desired result. When  $m = 1$ , since we are in the overparameterized regime,  $\mathbf{X}_p$  is a fat matrix. Thus, the regressor memorizes the training data and the training error is equal to zero.  $\mathbf{x}_0$  is part of training set, and so  $\hat{y}_0 = \mathbf{x}_0^\top \beta + \epsilon$ . Since  $\beta \sim \mathcal{N}(0, \frac{1}{D} \mathbf{I}_p)$ , we have that  $\mathbf{x}_0^\top \beta \sim \mathcal{N}(0, \frac{1}{D} \|\mathbf{x}_0\|^2)$ . Since  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ , we have that  $\hat{y}_0 = \mathbf{x}_0^\top \beta + \epsilon \sim \mathcal{N}(0, \frac{1}{D} \|\mathbf{x}_0\|^2 + \sigma^2)$ .

The probability distribution functions of the two Gaussians are then equal at  $\pm \alpha$ :

$$\begin{aligned}
\frac{1}{\sigma_0 \sqrt{2\pi}} \exp \left( -\frac{1}{2} \left( \frac{\alpha}{\sigma_0} \right)^2 \right) &= \frac{1}{\sigma_1 \sqrt{2\pi}} \exp \left( -\frac{1}{2} \left( \frac{\alpha}{\sigma_1} \right)^2 \right) \\
\frac{\sigma_1}{\sigma_0} &= \exp \left( -\frac{1}{2} \left( \left( \frac{\alpha}{\sigma_1} \right)^2 - \left( \frac{\alpha}{\sigma_0} \right)^2 \right) \right) \\
\frac{\sigma_1}{\sigma_0} &= \exp \left( -\frac{\alpha^2}{2} \frac{\sigma_0^2 - \sigma_1^2}{\sigma_0^2 \sigma_1^2} \right) \\
\log \left( \frac{\sigma_1}{\sigma_0} \right) &= -\frac{\alpha^2}{2} \frac{\sigma_0^2 - \sigma_1^2}{\sigma_0^2 \sigma_1^2} \\
\alpha &= \sqrt{\frac{2\sigma_0^2 \sigma_1^2 \log \left( \frac{\sigma_1}{\sigma_0} \right)}{\sigma_1^2 - \sigma_0^2}} \\
\alpha &= \sqrt{\frac{\sigma_0^2 \sigma_1^2 \log \left( \frac{\sigma_1^2}{\sigma_0^2} \right)}{\sigma_1^2 - \sigma_0^2}}.
\end{aligned}$$

The membership advantage is then derived by writing out the probabilities in Definition 2.1 in terms of the Gaussian cumulative distribution functions, noting that the decision region switches at  $\pm \alpha$ .  $\square$

*Proof of Lemma 4.4.* The lemma follows identically to Theorem 3.2 with an additional additive  $\bar{\sigma}^2$  to  $\sigma_0^2$  due to the noise added in the  $m = 0$  case. The remainder follows by plugging in  $p = D$  and applying Prop. 4.3 for the generalization error.  $\square$

### B.3 Proof of Theorem 4.1

*Proof of Theorem 4.1.* Let the input be  $\mathbf{x}_{0,p} \in \mathbb{R}^p$ . Similar to the proof of theorem 3.2, we can write

$$\mathbf{X}\beta + \epsilon = \mathbf{X}_p \beta_p + \mathbf{X}_{\bar{p}} \beta_{\bar{p}} + \epsilon = \mathbf{X}_p \beta_p + \eta$$

where  $\eta = \mathbf{X}_{\bar{p}}\beta_{\bar{p}} + \epsilon \sim \mathcal{N}(0, (1 + \sigma^2 - \frac{p}{D}) \mathbf{I}_n)$ . Hence, we have

$$\hat{y}_0 = \mathbf{x}_{0,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top (\mathbf{X}_p \beta_p + \eta) = \mathbf{x}_{0,p}^\top \mathbf{X}_p^\top \left( \mathbf{X}_p \mathbf{X}_p^\top + n\lambda \mathbf{I}_p \right)^{-1} (\mathbf{X}_p \beta_p + \eta). \quad (7)$$

First note that in the case  $m = 0$ , we have

$$\begin{aligned} \hat{y}_0 &= \mathbf{x}_{0,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top (\mathbf{X}_p \beta_p + \eta) \\ &= \mathbf{x}_{0,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \mathbf{X}_p \beta_p + \mathbf{x}_{0,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \eta. \end{aligned} \quad (8)$$

Letting the sample covariance matrix  $\hat{\Sigma} \triangleq (1/n) \mathbf{X}_p^\top \mathbf{X}_p$ , the first term in (8) can be written as

$$\begin{aligned} \mathbf{x}_{0,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \mathbf{X}_p \beta_p &= \mathbf{x}_{0,p}^\top \left( \hat{\Sigma} + \lambda \mathbf{I}_p \right)^{-1} \hat{\Sigma} \beta_p \\ &= \mathbf{x}_{0,p}^\top \left( \hat{\Sigma} + \lambda \mathbf{I}_p \right)^{-1} \left( \hat{\Sigma} + \lambda \mathbf{I}_p - \lambda \mathbf{I}_p \right) \beta_p \\ &= \left( \mathbf{x}_{0,p} - \lambda (\hat{\Sigma} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}_{0,p} \right)^\top \beta_p. \end{aligned}$$

Since  $\beta_p \sim (0, \frac{1}{D} \mathbf{I}_D)$  using Lemma B.2, this converges to a gaussian with zero mean and variance

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{D} \left\| \mathbf{x}_{0,p} - \lambda (\hat{\Sigma} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}_{0,p} \right\|_2^2 &= \lim_{n \rightarrow \infty} \frac{1}{D} \left\{ \left\| \mathbf{x}_{0,p} \right\|_2^2 - 2\lambda [\mathbf{x}_{0,p}^\top (\hat{\Sigma} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}_{0,p}] \right. \\ &\quad \left. + \lambda^2 [\mathbf{x}_{0,p}^\top (\hat{\Sigma} + \lambda \mathbf{I}_p)^{-2} \mathbf{x}_{0,p}] \right\} \\ &= \frac{\left\| \mathbf{x}_{0,p} \right\|_2^2}{D} (1 - 2\lambda g(-\lambda) + \lambda^2 g'(-\lambda)) \end{aligned}$$

where for the second equality, we have used the fact that using Lemma B.1 by setting  $\mathbf{C}_n = (1/n) \mathbf{x}_{0,p} \mathbf{x}_{0,p}^\top$ , as  $n \rightarrow \infty$ , almost surely,

$$\frac{1}{n} \mathbf{x}_{0,p}^\top (\hat{\Sigma} + \lambda \mathbf{I}_p)^{-1} \mathbf{x}_{0,p} \rightarrow \frac{1}{n} \left\| \mathbf{x}_{0,p} \right\|_2^2 g(-\lambda), \quad \frac{1}{n} \mathbf{x}_{0,p}^\top (\hat{\Sigma} + \lambda \mathbf{I}_p)^{-2} \mathbf{x}_{0,p} \rightarrow \frac{1}{n} \left\| \mathbf{x}_{0,p} \right\|_2^2 g'(-\lambda).$$

For the second term in (8), using the rotational invariance of gaussian distribution, without loss of generality, we can let  $\eta$  to be  $\mathbf{e}_1 \|\eta\|_2$ , where  $\mathbf{e}_1$  is the first standard unit vector. Now, note that we have

$$\|\eta\|_2 \mathbf{x}_{0,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \mathbf{e}_1 = \|\eta\|_2 \mathbf{x}_{0,p}^\top \left( \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top + \lambda n \mathbf{I}_p + \sum_{i=2}^n \mathbf{x}_{i,p} \mathbf{x}_{i,p}^\top \right)^{-1} \mathbf{x}_{1,p}$$

where  $\mathbf{x}_{i,p}^\top \in \mathbb{R}^p$  is the  $i$ 'th row of  $\mathbf{X}_p$ . Letting  $\mathbf{A}_\lambda \triangleq \lambda \mathbf{I}_p + \frac{1}{n} \sum_{i=2}^n \mathbf{x}_{i,p} \mathbf{x}_{i,p}^\top$ , by using the Sherman-Morrison formula, we have

$$\begin{aligned} \|\eta\|_2 \mathbf{x}_{0,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \mathbf{e}_1 &= \|\eta\|_2 \mathbf{x}_{0,p}^\top (\mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top + n \mathbf{A}_\lambda)^{-1} \mathbf{x}_{1,p} \\ &= \frac{\|\eta\|_2}{n} \mathbf{x}_{0,p}^\top \left( \mathbf{A}_\lambda^{-1} - \frac{\mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \right) \mathbf{x}_{1,p} \\ &= \frac{\|\eta\|_2}{n} \mathbf{x}_{0,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \left( 1 - \frac{\mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \right) \\ &= \|\eta\|_2 \frac{\mathbf{x}_{0,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}}. \end{aligned}$$



Note that using Lemma B.1 by setting  $\mathbf{C}_n = (1/n)\mathbf{x}_{1,p}\mathbf{x}_{1,p}^\top$  and  $\mathbf{C}_n = (1/n)\mathbf{x}_{0,p}\mathbf{x}_{0,p}^\top$ , respectively, for  $n, p \rightarrow \infty$ , almost surely,

$$\frac{1}{n}\mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \rightarrow \gamma g(-\lambda), \quad \frac{1}{n}\mathbf{x}_{0,p}^\top \mathbf{A}_\lambda^{-2} \mathbf{x}_{0,p} \rightarrow \frac{\|\mathbf{x}_{0,p}\|_2^2}{n} g'(-\lambda).$$

Thus, since  $\mathbf{x}_{1,p} \sim (0, \mathbf{I}_p)$ , using Lemma B.2,  $\|\eta\|_2 \mathbf{x}_{0,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \mathbf{e}_1$  converges to a gaussian with mean zero and variance

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{\|\eta\|_2^2}{n^2(1 + \gamma g(-\lambda))^2} \|\mathbf{A}_\lambda^{-1} \mathbf{x}_{0,p}\|_2^2 &= \frac{\|\eta\|_2^2 \|\mathbf{x}_{0,p}\|_2^2}{n^2(1 + \gamma g(-\lambda))^2} g'(-\lambda) \\ &= \frac{\|\mathbf{x}_{0,p}\|_2^2}{p} \frac{g'(-\lambda)\gamma}{1 + \gamma g(-\lambda))^2} \left( \sigma^2 + 1 - \frac{p}{D} \right). \end{aligned}$$

Hence, by independence of  $\beta_p$  and  $\eta$ , for  $m = 0$ , as  $n \rightarrow \infty$ , such that  $p/n = \gamma$ , the output  $\hat{y}_0$  as in (7), converges in distribution to a gaussian with mean zero and variance

$$\frac{g'(-\lambda)\gamma}{(1 + g(-\lambda)\gamma)^2} \left( \sigma^2 + 1 - \frac{p}{D} \right) \frac{\|\mathbf{x}_{0,p}\|_2^2}{p} + (1 - 2\lambda g(-\lambda) + \lambda^2 g'(-\lambda)) \frac{\|\mathbf{x}_{0,p}\|_2^2}{D}.$$

Now consider the  $m = 1$  case where the input belongs to training data. Without loss of generality, assume that the input is the first row of  $\mathbf{X}_p$ , i.e.  $\mathbf{x}_0 := \mathbf{x}_1$ . Note that in this case for  $\eta = \mathbf{X}_{\bar{p}}\boldsymbol{\beta}_{\bar{p}} + \epsilon$ , we have  $\eta_i \sim \mathcal{N}\left(0, \left(\sigma^2 + \frac{\|\mathbf{x}_{1,\bar{p}}\|_2^2}{D}\right) \mathbf{I}_n\right)$ , for  $i = 1$ ,  $\eta_i \sim \mathcal{N}\left(0, \left(1 + \sigma^2 - \frac{p}{D}\right)\right)$ , for  $i = 2, 3, \dots, n$ , and  $\eta_i$ 's are independent. We have

$$\begin{aligned} \hat{y}_0 &= \mathbf{x}_{1,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top (\mathbf{X}_p \boldsymbol{\beta}_p + \eta) \\ &= \mathbf{x}_{1,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \mathbf{X}_p \boldsymbol{\beta}_p + \mathbf{x}_{1,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \eta \\ &= \mathbf{x}_{1,p}^\top \left( \hat{\boldsymbol{\Sigma}} + \lambda \mathbf{I}_p \right)^{-1} \hat{\boldsymbol{\Sigma}} \boldsymbol{\beta}_p + \mathbf{x}_{1,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \eta. \end{aligned} \tag{9}$$

The first term in (9) can be written as

$$\begin{aligned} \mathbf{x}_{1,p}^\top \left( \hat{\boldsymbol{\Sigma}} + \lambda \mathbf{I}_p \right)^{-1} \hat{\boldsymbol{\Sigma}} \boldsymbol{\beta}_p &= \mathbf{x}_{1,p}^\top \boldsymbol{\beta}_p - \lambda \mathbf{x}_{1,p}^\top \left( \hat{\boldsymbol{\Sigma}} + \lambda \mathbf{I}_p \right)^{-1} \boldsymbol{\beta}_p \\ &= \mathbf{x}_{1,p}^\top \boldsymbol{\beta}_p - \lambda \mathbf{x}_{1,p}^\top \left[ \frac{1}{n} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top + \mathbf{A}_\lambda \right]^{-1} \boldsymbol{\beta}_p \\ &= \mathbf{x}_{1,p}^\top \boldsymbol{\beta}_p - \lambda \mathbf{x}_{1,p}^\top \left[ \mathbf{A}_\lambda^{-1} - \frac{\mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \right] \boldsymbol{\beta}_p \\ &= \mathbf{x}_{1,p}^\top (\mathbf{I}_p - \lambda \mathbf{A}_\lambda^{-1}) \boldsymbol{\beta}_p + \frac{\lambda \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \boldsymbol{\beta}_p}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \\ &= \mathbf{x}_{1,p}^\top \left( \mathbf{I}_p - \frac{\lambda}{1 + (1/n) \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \mathbf{A}_\lambda^{-1} \right) \boldsymbol{\beta}_p \end{aligned}$$

Hence, since  $\boldsymbol{\beta}_i \sim (0, \frac{1}{D} \mathbf{I}_D)$ , using Lemma B.2, The first term in (9) converges to a gaussian with zero mean and variance

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{D} \left\| \hat{\boldsymbol{\Sigma}} \left( \hat{\boldsymbol{\Sigma}} + \lambda \mathbf{I}_p \right)^{-1} \mathbf{x}_{1,p} \right\|_2^2 &= \lim_{n \rightarrow \infty} \frac{1}{D} \left\| \left( \mathbf{I}_p - \frac{\lambda}{1 + (1/n) \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \mathbf{A}_\lambda^{-1} \right) \mathbf{x}_{1,p} \right\|_2^2 \\ &= \frac{\|\mathbf{x}_{1,p}\|_2^2}{D} \left( 1 - \frac{2\lambda g(-\lambda)}{1 + \frac{\gamma \|\mathbf{x}_{1,p}\|_2^2}{p} g(-\lambda)} + \frac{\lambda^2 g'(-\lambda)}{\left( 1 + \frac{\gamma \|\mathbf{x}_{1,p}\|_2^2}{p} g(-\lambda) \right)^2} \right). \end{aligned}$$

where for the second equality we have used the fact that using Lemma B.1 by setting  $\mathbf{C}_n = (1/n)\mathbf{x}_{1,p}\mathbf{x}_{1,p}^\top$ , as  $n \rightarrow \infty$ , such that  $p/n = \gamma$ , almost surely,

$$\frac{1}{n}\mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \rightarrow \frac{1}{n}\|\mathbf{x}_{1,p}\|_2^2 g(-\lambda), \quad \frac{1}{n}\mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-2} \mathbf{x}_{1,p} \rightarrow \frac{1}{n}\|\mathbf{x}_{1,p}\|_2^2 g'(-\lambda).$$

Now consider the second term in (9). It can be written as

$$\begin{aligned} \mathbf{x}_{1,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{X}_p^\top \boldsymbol{\eta} &= \mathbf{x}_{1,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{x}_{1,p}^\top \eta_1 \\ &\quad + \sum_{i=2}^n \mathbf{x}_{1,p}^\top \left( \mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p \right)^{-1} \mathbf{x}_{i,p}^\top \eta_i \\ &= A\eta_1 + \sum_{i=2}^n B_i \eta_i. \end{aligned}$$

First consider

$$\begin{aligned} A &= \frac{1}{n}\mathbf{x}_{1,p}^\top \left( \frac{1}{n}\mathbf{x}_{1,p}\mathbf{x}_{1,p}^\top + \lambda \mathbf{I}_p + \underbrace{\frac{1}{n} \sum_{i=2}^n \mathbf{x}_{i,p}\mathbf{x}_{i,p}^\top}_{\mathbf{A}_\lambda} \right)^{-1} \mathbf{x}_{i,p} \\ &= \frac{1}{n}\mathbf{x}_{1,p}^\top \left( \mathbf{A}_\lambda^{-1} - \frac{\mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}\mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \right) \mathbf{x}_{1,p} \\ &= \frac{1}{n}\mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p} \left( 1 - \frac{\mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \right) \\ &= \frac{\mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}}{n + \mathbf{x}_{1,p}^\top \mathbf{A}_\lambda^{-1} \mathbf{x}_{1,p}} \xrightarrow{a.s.} \frac{\gamma g(-\lambda)(\|\mathbf{x}_{1,p}\|_2^2/p)}{1 + \gamma g(-\lambda)(\|\mathbf{x}_{1,p}\|_2^2/p)}. \end{aligned}$$

Now, consider

$$\begin{aligned} B_2 &= \frac{1}{n}\mathbf{x}_{1,p}^\top \left( \frac{1}{n}\mathbf{U}\mathbf{U}^\top + \lambda \mathbf{I}_p + \underbrace{\frac{1}{n} \sum_{i=3}^n \mathbf{x}_{i,p}\mathbf{x}_{i,p}^\top}_{\mathbf{A}_{2,\lambda}} \right)^{-1} \mathbf{x}_{i,p}; \quad \mathbf{U} \triangleq [\mathbf{x}_{1,p} \ \mathbf{x}_{2,p}] \\ &= \frac{1}{n}\mathbf{x}_{1,p}^\top \left[ \mathbf{A}_{2,\lambda}^{-1} - \underbrace{\mathbf{A}_{2,\lambda}^{-1} \mathbf{U} \left( n\lambda \mathbf{I}_p + \mathbf{U}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{U} \right)^{-1} \mathbf{U}^\top \mathbf{A}_{2,\lambda}^{-1}}_{\mathbf{C}_2} \right] \mathbf{x}_{2,p}. \end{aligned}$$

We have

$$\begin{aligned} \mathbf{C}_2 &= [\mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \ \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}] \begin{bmatrix} n\lambda + \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} & \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \\ \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} & n\lambda + \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \\ \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \end{bmatrix} \\ &= \left[ \underbrace{\left( n\lambda + \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \right) \left( n\lambda + \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \right) - \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p}}_{D_2} \right]^{-1} \tilde{\mathbf{C}}_2. \end{aligned}$$

We have

$$\begin{aligned}
\tilde{\mathbf{C}}_2 &= [\mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \quad \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}] \begin{bmatrix} n\lambda + \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} & -\mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \\ -\mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} & n\lambda + \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \end{bmatrix} \begin{bmatrix} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \\ \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \end{bmatrix} \\
&= \left( n\lambda + \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \right) \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} - \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \\
&\quad - \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} + \left( n\lambda + \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \right) \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1}.
\end{aligned}$$

Thus,

$$\begin{aligned}
B_2 &= \frac{1}{n} \mathbf{x}_{1,p}^\top \left[ \mathbf{A}_{2,\lambda}^{-1} - \frac{\tilde{\mathbf{C}}_2}{D_2} \right] \mathbf{x}_{2,p} \\
&= \frac{\mathbf{x}_{1,p}^\top}{n} \left\{ \mathbf{A}_{2,\lambda}^{-1} - \left[ \left( n\lambda + \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \right) \left( n\lambda + \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \right) - \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \right]^{-1} \right. \\
&\quad \left[ \left( n\lambda + \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \right) \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} - \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \right. \\
&\quad \left. \left. - \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} + \left( n\lambda + \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p} \right) \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \mathbf{x}_{2,p}^\top \mathbf{A}_{2,\lambda}^{-1} \right] \right\} \mathbf{x}_{2,p}.
\end{aligned}$$

using Lemma B.1 by setting  $\mathbf{C}_n = (1/n) \mathbf{x} \mathbf{x}^\top$ , as  $n, p \rightarrow \infty$ , such that  $p/n = \gamma$ , almost surely,

$$\frac{1}{n} \mathbf{x}^\top \mathbf{A}_\lambda^{-2} \mathbf{x} \rightarrow \frac{\|\mathbf{x}\|_2^2}{n} g(-\lambda).$$

Hence, letting  $n \rightarrow \infty$ ,  $B_2$  converges weakly to

$$\begin{aligned}
B'_2 &= \frac{1}{n} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \left\{ \mathbf{A}_{2,\lambda} - \left[ \left( \lambda + g(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{n} \right) \left( \lambda + g(-\lambda) \frac{\|\mathbf{x}_{2,p}\|_2^2}{n} \right) - \left( \frac{\mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}}{n} \right)^2 \right]^{-1} \right. \\
&\quad \left[ \left( \lambda + g(-\lambda) \frac{\|\mathbf{x}_{2,p}\|_2^2}{n} \right) \frac{\mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top}{n} - \left( \frac{\mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}}{n} \right) \left( \frac{\mathbf{x}_{1,p} \mathbf{x}_{2,p}^\top}{n} + \frac{\mathbf{x}_{2,p} \mathbf{x}_{1,p}^\top}{n} \right) \right. \\
&\quad \left. \left. + \left( \lambda + g(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{n} \right) \frac{\mathbf{x}_{2,p} \mathbf{x}_{2,p}^\top}{n} \right] \right\} \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \\
&= \frac{1}{n} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \frac{\lambda^2}{\left( \lambda + g(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{n} \right) \left( \lambda + g(-\lambda) \frac{\|\mathbf{x}_{2,p}\|_2^2}{n} \right) - \left( \frac{\mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}}{n} \right)^2}.
\end{aligned}$$

Note that by LLN,  $(1/n) \|\mathbf{x}_{2,p}\|_2^2 \rightarrow \gamma$  and  $(1/n^2) \left( \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \right)^2 \rightarrow 0$ . Further, since  $\mathbf{x}_{2,p} \sim (0, \mathbf{I}_p)$ , using Lemma B.2, as  $n \rightarrow \infty$ ,  $(1/\sqrt{n}) \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p}$  converges to a gaussian with mean zero and variance

$$\lim_{n \rightarrow \infty} \frac{1}{n} \|\mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{1,p}\|_2^2 = \lim_{n \rightarrow \infty} \frac{1}{n} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-2} \mathbf{x}_{1,p} = \frac{1}{n} \|\mathbf{x}_{1,p}\|_2^2 g'(-\lambda)$$

where for the second equality we have used Lemma B.1 with  $\mathbf{C}_n = (1/n) \mathbf{x}_{1,p} \mathbf{x}_{1,p}^\top$ . Hence,  $B_2$  converges weakly to

$$\frac{\tilde{\sigma}}{n} \mathbf{x}_{1,p}^\top \mathbf{A}_{2,\lambda}^{-1} \mathbf{x}_{2,p} \xrightarrow{\mathcal{D}} \mathcal{N} \left( 0, \frac{\gamma \tilde{\sigma}^2 \|\mathbf{x}_{1,p}\|_2^2}{pn} g'(-\lambda) \right)$$

where

$$\tilde{\sigma} = \frac{\lambda^2}{\lambda^2 + \gamma \lambda g(-\lambda) \left[ 1 + \frac{\|\mathbf{x}_{1,p}\|_2^2}{p} \right] + \gamma^2 m^2(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{p}}.$$

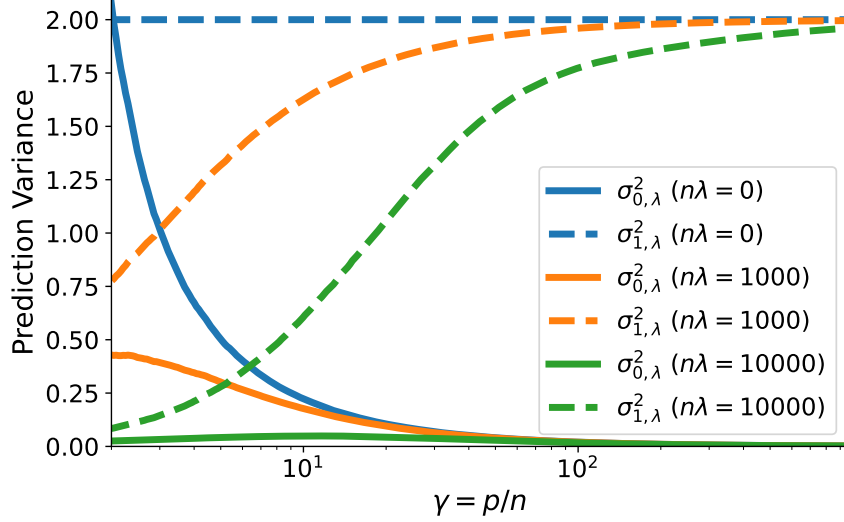


Figure 6: Theoretical variances of the predictions  $\hat{y}_0$  by ridge regularized linear regression models for the Gaussian data setting with  $n = 10^3$ ,  $D = 10^7$ , and  $\sigma = 1$  on a single sampled  $\mathbf{x}_0$  for when  $\mathbf{x}_0$  is a test point ( $\sigma_{0,\lambda}^2$ ) and when  $\mathbf{x}_0$  is a training point ( $\sigma_{1,\lambda}^2$ ) for different amounts of regularization  $\lambda$ . While increased ridge regularization decreases the variance  $\sigma_{0,\lambda}^2$  on training point predictions, it also decreases the variance  $\sigma_{1,\lambda}^2$  for test points in such a way that the two distributions become easier to distinguish. As such, membership inference is easier for ridge regularized models in this setting.

By symmetry over  $i$ ,  $\sum_{i=2}^n B_i^2 = (n-1)B_2^2$  that converges almost surely to  $\frac{\gamma\tilde{\sigma}^2\|\mathbf{x}_{1,p}\|_2^2}{p}g'(-\lambda)$ . Therefore using Lemma B.2

$$\sum_{i=2}^n B_i \eta_i \xrightarrow{\mathcal{D}} \mathcal{N}\left(0, \left(\sigma^2 + 1 - \frac{p}{D}\right) \gamma \tilde{\sigma}^2 \frac{\|\mathbf{x}_{1,p}\|_2^2 g'(-\lambda)}{p}\right).$$

Thus, by independence of  $\beta_p$  and  $\eta_i$ 's, as  $n, p \rightarrow \infty$  the output  $y$  converges in distribution to a gaussian with zero mean and variance

$$\begin{aligned} & \left[ \left( \frac{\lambda^2}{(\lambda + \gamma g(-\lambda))(\lambda + \gamma \frac{\|\mathbf{x}_{1,p}\|_2^2}{p} g(-\lambda))} \right)^2 \gamma g'(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{p} \right] \left( \sigma^2 + 1 - \frac{p}{D} \right) \\ & + \left( \frac{\gamma g(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{p}}{1 + \gamma g(-\lambda) \frac{\|\mathbf{x}_{1,p}\|_2^2}{p}} \right)^2 \left( \sigma^2 + \frac{\|\mathbf{x}_{1,p}\|_2^2}{D} \right) \\ & + \left( 1 - \frac{2\lambda g(-\lambda)}{1 + \gamma \frac{\|\mathbf{x}_{1,p}\|_2^2}{p} g(-\lambda)} + \frac{\lambda^2 g'(-\lambda)}{\left(1 + \gamma \frac{\|\mathbf{x}_{1,p}\|_2^2}{p} g(-\lambda)\right)^2} \right) \frac{\|\mathbf{x}_{1,p}\|_2^2}{D}, \end{aligned}$$

which completes the proof.  $\square$

## C Posterior Distributions in Non-Asymptotic Regime

Let  $f_{a|b}$  denote the probability density function of a random variable  $a$  conditioned on  $b$ . The following lemma derives the non-asymptotic probability densities of the prediction output of minimum norm least squares, conditioned on the  $m = 0$  and  $m = 1$  events and the choice of test point  $\mathbf{x}_0$ . For a matrix  $\mathbf{X} \in \mathbb{R}^{n \times D}$  and  $p \leq D$ , let  $\mathbf{X}_p$  denote the submatrix of the first  $p$  columns of  $\mathbf{X}$ . For a vector  $\mathbf{x} \in \mathbb{R}^D$ , let  $\mathbf{x}_p \in \mathbb{R}^p$  be defined accordingly.

**Lemma C.1.** *Let  $\hat{\beta}$  denote the minimum norm least squares interpolator computed from a random design matrix  $\mathbf{X} \in \mathbb{R}^{n \times D}$  and data  $\mathbf{y}$ . Conditioned on  $n < p \leq D$  and on  $\mathbf{x}_0$ , we have that  $\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{x}_0, \{m = 1\} \sim \mathcal{N}(0, \sigma_1^2)$ , where  $\sigma_1$  is defined as in Theorem 3.2. Furthermore,*

$$\begin{aligned} & f_{\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{x}_0, \{m=0\}}(x) \\ &= D^{\frac{D}{2}} \int_{\mathbb{R}^{n \times D}} \int_{\mathbb{R}^D} \frac{\exp \left[ -\frac{1}{2} \left[ \left( \frac{x - \mathbf{x}_{0,p}^\top \mathbf{X}_p^\top (\mathbf{X}_p \mathbf{X}_p^\top)^{-1} \mathbf{X} \beta}{\sigma \|\mathbf{x}_{0,p}\| \mathbf{X}_p (\mathbf{X}_p \mathbf{X}_p^\top)^{-2} \mathbf{X}_p} \right)^2 + D \beta^\top \beta + \text{Tr}(\mathbf{X}^\top \mathbf{X}) \right] \right]}{\sigma(2\pi)^{\frac{nD+D+1}{2}} \|\mathbf{x}_{0,p}\| \mathbf{X}_p (\mathbf{X}_p \mathbf{X}_p^\top)^{-2} \mathbf{X}_p} d\beta d\mathbf{X}. \end{aligned}$$

*Remark C.2.* While the density of  $\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{x}_0, \{m = 0\}$  cannot be written in a closed form, one may easily sample according to it, by first, sampling random  $\mathbf{X}$ ,  $\beta$  and then computing the minimum norm least squares interpolator.

*Proof of Lemma C.1.* Recall that conditioned on the design matrix  $\mathbf{X}$  and true coefficients  $\beta$ , the labels  $\mathbf{y}$  follow  $\mathbf{y} \mid \mathbf{X}, \beta \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$ . Then, the minimum norm least squares solution  $\hat{\beta}$  using the first  $p$  features follows

$$\hat{\beta} \mid \mathbf{X}, \beta \sim \mathcal{N} \left( \mathbf{X}_p^\top (\mathbf{X}_p \mathbf{X}_p^\top)^{-1} \mathbf{X} \beta, \sigma^2 \mathbf{X}_p^\top (\mathbf{X}_p \mathbf{X}_p^\top)^{-2} \mathbf{X}_p \right).$$

Hence, for the  $m = 0$  case where a fresh point  $\mathbf{x}_0$  is sampled, we have that the distribution of the model output conditioned on the design matrix  $\mathbf{X}$  and true coefficients  $\beta$  is

$$\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m = 0\} \sim \mathcal{N} \left( \mathbf{x}_{0,p}^\top \mathbf{X}_p^\top (\mathbf{X}_p \mathbf{X}_p^\top)^{-1} \mathbf{X} \beta, \sigma^2 \|\mathbf{x}_{0,p}\|^2 \mathbf{X}_p^\top (\mathbf{X}_p \mathbf{X}_p^\top)^{-2} \mathbf{X}_p \right)$$

for  $\|\mathbf{x}\|_{\mathbf{A}} := \sqrt{\mathbf{x}^\top \mathbf{A} \mathbf{x}}$  for any semidefinite matrix  $\mathbf{A}$  where we have additionally conditioned over any randomness in the choice of  $\mathbf{x}_0$ .

In the  $m = 1$  case, where  $\mathbf{x}_0$  is sampled uniformly from the rows of  $\mathbf{X}$ , we have that  $\mathbf{x}_{0,p}^\top \hat{\beta} = y_0 = \mathbf{x}_0^\top \beta + \epsilon$ , the associated label for  $\mathbf{x}_0$  since the linear regressor interpolates the training data. Hence

$$\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m = 1\} \equiv \mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{x}_0, \{m = 1\} \sim \mathcal{N} \left( 0, \frac{\|\mathbf{x}_0\|^2}{D} + \sigma^2 \right)$$

Let  $f_{\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m=0\}}$  denote the pdf of the random variable  $\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m = 0\}$  and  $f_{\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m=1\}}$  be defined similarly. Let  $f_{\mathbf{X}}$  denote the density of  $\mathbf{X}$ , a standard matrix-normal random variable, and let  $f_{\beta}$  denote the density of  $\beta \sim \mathcal{N}(0, \frac{1}{D} \mathbf{I}_D)$ . Then, we have that

$$\begin{aligned} & f_{\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{x}_0, \{m=0\}}(x) \\ &= \int_{\mathbb{R}^{n \times D}} \int_{\mathbb{R}^D} f_{\mathbf{x}_{0,p}^\top \hat{\beta} \mid \mathbf{X}, \beta, \mathbf{x}_0, \{m=0\}} f_{\beta} f_{\mathbf{X}} d\beta d\mathbf{X} \\ &= D^{\frac{D}{2}} \int_{\mathbb{R}^{n \times D}} \int_{\mathbb{R}^D} \frac{\exp \left[ -\frac{1}{2} \left[ \left( \frac{x - \mathbf{x}_{0,p}^\top \mathbf{X}_p^\top (\mathbf{X}_p \mathbf{X}_p^\top)^{-1} \mathbf{X} \beta}{\sigma \|\mathbf{x}_{0,p}\| \mathbf{X}_p (\mathbf{X}_p \mathbf{X}_p^\top)^{-2} \mathbf{X}_p} \right)^2 + D \beta^\top \beta + \text{Tr}(\mathbf{X}^\top \mathbf{X}) \right] \right]}{\sigma(2\pi)^{\frac{nD+D+1}{2}} \|\mathbf{x}_{0,p}\| \mathbf{X}_p (\mathbf{X}_p \mathbf{X}_p^\top)^{-2} \mathbf{X}_p} d\beta d\mathbf{X}. \end{aligned}$$

□

**Lemma C.3.** Let  $\hat{\beta}_\lambda = (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{y}$  denote ridge regularized least squares estimator computed from random design matrix  $\mathbf{X} \in \mathbb{R}^{n \times D}$ , data  $\mathbf{y}$ , and subset of first  $p$  features. Conditioned on the choice of test point  $\mathbf{x}_0$ , we have that in the  $m = 0$  case, where a fresh test point is drawn from the data distribution,

$$f_{\mathbf{x}_{0,p}^\top \hat{\beta} | \mathbf{x}_0, \{m=0\}}(x) = \frac{D^{\frac{D}{2}}}{\sigma(2\pi)^{\frac{nD+D+1}{2}}} \times \int_{\mathbb{R}^{n \times D}} \int_{\mathbb{R}^D} \frac{\exp \left[ -\frac{1}{2} \left[ \left( \frac{x - \mathbf{x}_{0,p}^\top (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X} \beta}{\sigma \|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1}} \right)^2 + D\beta^\top \beta + \text{Tr}(\mathbf{X}^\top \mathbf{X}) \right] \right]}{\|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1}} d\beta d\mathbf{X}.$$

Furthermore, conditioned on  $m = 1$  when  $\mathbf{x}_0$  is a row of  $\mathbf{X}$  we have that

$$f_{\mathbf{x}_{0,p}^\top \hat{\beta} | \mathbf{x}_0, \{m=1\}}(x) = \frac{D^{\frac{D}{2}}}{\sigma(2\pi)^{\frac{nD+1}{2}}} \times \int_{\mathbb{R}^{(n-1) \times D}} \int_{\mathbb{R}^D} \frac{\exp \left[ -\frac{1}{2} \left[ \left( \frac{x - \mathbf{x}_{0,p}^\top (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X} \beta}{\sigma \|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1}} \right)^2 + D\beta^\top \beta + \text{Tr}(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) \right] \right]}{\|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I})^{-1}} d\beta d\tilde{\mathbf{X}}.$$

where without loss of generality, we take  $\mathbf{x}_0$  to be the first row of  $\mathbf{X}$  and  $\tilde{\mathbf{X}}$  to denote the matrix of the final  $n - 1$  rows of  $\mathbf{X}$ .

*Remark C.4.* As in the case of Lemma C.1, one can efficiently sample from the above distribution by first drawing a Gaussian random matrix  $\mathbf{X}$ , the Gaussian random vector  $\beta$ , the Bernoulli random variable  $m$ , and then either a new test point  $\mathbf{x}_0$  or a row of  $\mathbf{X}$  and learning the ridge-regularized estimator  $\hat{\beta}_\lambda$ .

*Proof of Lemma C.3.* Note that conditioned on the design matrix  $\mathbf{X}$  and the true coefficients  $\beta$ , the labels  $\mathbf{y}$  follow  $\mathbf{y} | \mathbf{X}, \beta \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$ . Next, for

$$\hat{\beta}_\lambda := (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \mathbf{X}_p^\top \mathbf{y}$$

the  $\lambda$ -ridge regularized estimator, we have that

$$\hat{\beta}_\lambda | \mathbf{X}, \beta \sim \mathcal{N} \left( (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \mathbf{X}_p^\top \mathbf{X} \beta, \sigma^2 (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \right).$$

Hence,

$$\begin{aligned} & \mathbf{x}_{0,p}^\top \hat{\beta}_\lambda | \mathbf{X}, \mathbf{x}_0, \beta \\ & \sim \mathcal{N} \left( \mathbf{x}_{0,p}^\top (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \mathbf{X}_p^\top \mathbf{X} \beta, \sigma^2 \mathbf{x}_{0,p}^\top (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_p + n\lambda \mathbf{I}_p)^{-1} \mathbf{x}_{0,p} \right), \end{aligned}$$

where we have additionally conditioned over any randomness in the choice of  $\mathbf{x}_0$ .

In the  $m = 0$  case, where  $\mathbf{x}_0$  is a freshly drawn point independent of the data  $\mathbf{X}$ , we may marginalize to remove the conditioning. Let  $f_{\mathbf{x}_{0,p}^\top \hat{\beta}_\lambda | \mathbf{X}, \beta, \mathbf{x}_0, \{m=0\}}$  denote the probability density function of the random variable  $\mathbf{x}_{0,p}^\top \hat{\beta}_\lambda | \mathbf{X}, \beta, \mathbf{x}_0, \{m=0\}$  and  $f_{\mathbf{x}_{0,p}^\top \hat{\beta} | \mathbf{X}, \mathbf{x}_0, \{m=1\}}$  be defined similarly. Let  $f_{\mathbf{X}}$  denote the density of  $\mathbf{X}$ , a standard matrix-normal random variable, and let  $f_\beta$  denote the density of  $\beta \sim \mathcal{N}(0, \frac{1}{D} \mathbf{I}_D)$ . Then we have that

$$f_{\mathbf{x}_{0,p}^\top \hat{\beta} | \mathbf{x}_0, \{m=0\}}(x) = \int_{\mathbb{R}^{n \times D}} \int_{\mathbb{R}^D} f_{\mathbf{x}_{0,p}^\top \hat{\beta} | \mathbf{X}, \mathbf{x}_0, \{m=0\}} f_\beta f_{\mathbf{X}} d\beta d\mathbf{X}.$$

Thus,

$$f_{\mathbf{x}_{0,p}^\top \hat{\beta} | \mathbf{x}_0, \{m=0\}}(x) = \frac{D^{\frac{D}{2}}}{\sigma(2\pi)^{\frac{nD+D+1}{2}}} \times \int_{\mathbb{R}^{n \times D}} \int_{\mathbb{R}^D} \frac{\exp \left[ -\frac{1}{2} \left[ \left( \frac{x - \mathbf{x}_{0,p}^\top (\mathbf{X}_p^\top \mathbf{X}_{p+n\lambda I})^{-1} \mathbf{X}_p^\top \mathbf{X} \beta}{\sigma \|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_{p+n\lambda I})^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_{p+n\lambda I})^{-1}} \right)^2 + D\beta^\top \beta + \text{Tr}(\mathbf{X}^\top \mathbf{X}) \right] \right]}{\|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_{p+n\lambda I})^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_{p+n\lambda I})^{-1}} d\beta d\mathbf{X}.$$

In the  $m = 1$  case, because  $\mathbf{x}_0$  is a row of  $\mathbf{X}$ , we condition on  $\mathbf{x}_0$  but not on the remaining rows of  $\mathbf{X}$ . Without loss of generality, let  $\mathbf{x}_0$  be the first row of  $\mathbf{X}$  which can be done since  $\mathbf{x}_0$  is selected uniformly and the rows of  $\mathbf{X}$  are independent and identically distributed. Let  $\tilde{\mathbf{X}} \in \mathbb{R}^{(n-1) \times D}$  denote  $\mathbf{X}$  with its first row omitted such that  $\mathbf{X} = [\mathbf{x}_0; \tilde{\mathbf{X}}]$ . Following the same approach as the preceding marginalization, we have that

$$f_{\mathbf{x}_{0,p}^\top \hat{\beta} | \mathbf{x}_0, \{m=1\}}(x) = \int_{\mathbb{R}^{(n-1) \times D}} \int_{\mathbb{R}^D} f_{\mathbf{x}_{0,p}^\top \hat{\beta} | \mathbf{X}, \mathbf{x}_0, \{m=1\}} f_{\tilde{\mathbf{X}}} d\beta d\tilde{\mathbf{X}}$$

Thus,

$$f_{\mathbf{x}_{0,p}^\top \hat{\beta} | \mathbf{x}_0, \{m=1\}}(x) = \frac{D^{\frac{D}{2}}}{\sigma(2\pi)^{\frac{nD+1}{2}}} \times \int_{\mathbb{R}^{(n-1) \times D}} \int_{\mathbb{R}^D} \frac{\exp \left[ -\frac{1}{2} \left[ \left( \frac{x - \mathbf{x}_{0,p}^\top (\mathbf{X}_p^\top \mathbf{X}_{p+n\lambda I})^{-1} \mathbf{X}_p^\top \mathbf{X} \beta}{\sigma \|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_{p+n\lambda I})^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_{p+n\lambda I})^{-1}} \right)^2 + D\beta^\top \beta + \text{Tr}(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}) \right] \right]}{\|\mathbf{x}_{0,p}\| (\mathbf{X}_p^\top \mathbf{X}_{p+n\lambda I})^{-1} \mathbf{X}_p^\top \mathbf{X}_p (\mathbf{X}_p^\top \mathbf{X}_{p+n\lambda I})^{-1}} d\beta d\tilde{\mathbf{X}}.$$

□

## D Experimental Implementation Details

All experiments were ran only on CPUs on our internal servers without GPU processing. Processors used may have included Intel Xeon CPU E5-2630 (256GB RAM), Intel Xeon Silver 4214 CPU (192GB RAM), Intel Xeon Platinum 8260 CPU (192GB RAM), and AMD Ryzen Threadripper 1900X (32GB RAM). Our code is primarily written in Python and mainly uses numpy implementations of linear algebra operations. Please refer to our code on the Github page for more details.

The histograms in Figure 1 are obtained as follows. We first sample a vector  $\mathbf{x}_0 \sim \mathcal{N}(0, \mathbf{I}_D)$ , where  $D = 20,000$ . Then, for each  $p = \gamma n$ , we perform the following procedure 20,000 times. We sample  $\beta \sim \mathcal{N}(0, \frac{1}{D} \mathbf{I}_D)$ . Then, we sample an  $n \times D$  matrix  $\mathbf{X}$  such that each element is iid standard normal. We then generate the ground truth vector  $\mathbf{y} = \mathbf{X}\beta + \epsilon$ , where  $\epsilon$  is an  $n$ -dimensional vector whose elements are iid standard normal. We obtain the least squares estimates  $\hat{\beta}$  on the first  $p$  columns of  $\mathbf{X}$  and on the vector  $\mathbf{y}$  using numpy's `lstsq` function. Finally, we collect the  $\hat{y}_0 = \mathbf{x}_{0,p}^\top \hat{\beta}$  of all 20,000 models to form the blue histograms in Figure 1. The orange histograms are formed the same way except that the first row of  $\mathbf{X}$  is replaced with  $\mathbf{x}_0$  and the first element of  $\mathbf{y}$  is replaced with  $y_0 = \mathbf{x}_0^\top \beta + \epsilon_0$  for  $\epsilon_0 \sim \mathcal{N}(0, 1)$ .

The experiment in Figure 2b is performed as follows. In the experiment, we estimate the optimal membership advantage. Since the optimal MI adversary requires knowledge of the linear regression model's output distributions when a data point  $\mathbf{x}_0$  is in its training dataset ( $m = 1$ ) and when  $\mathbf{x}_0$  is not ( $m = 0$ ), we approximate these distributions by forming discrete histograms. To obtain the samples for the histograms, we use the same procedure as detailed in the previous paragraph, except that we obtain 100,000 samples for each histogram for increased precision. From these samples, the discrete histograms for  $(\hat{y}_0 | m = 0)$  and  $(\hat{y}_0 | m = 1)$  for a given  $\gamma$  are then formed by splitting the interval between the minimum and maximum values over both  $(\hat{y}_0 | m = 0)$  and  $(\hat{y}_0 | m = 1)$  into 150 equally spaced bins. The histograms are normalized so that they represent probability mass functions (i.e. the bin counts sum to 1). Finally, treating the two

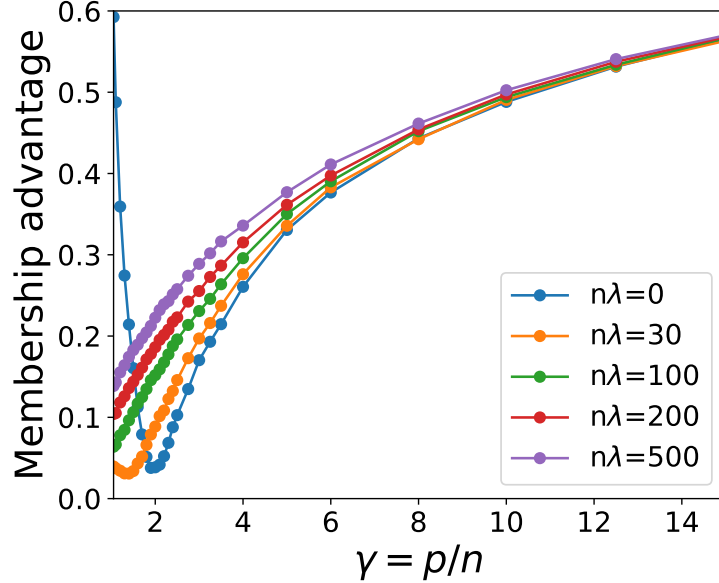


Figure 7: Experimental membership advantages for ridge-regularized linear regression on Gaussian data with  $n = 100$ ,  $D = 3000$ , and  $\sigma = 1$  for different regularization strengths  $\lambda$ . As predicted by our theory, membership advantage increases with additional regularization in the sufficiently overparameterized regime. This experiment verifies our theoretical findings.

histograms as probability mass functions, the membership advantage is calculated according to Definition 2.1. For Figure 2b, this procedure is repeated 20 times, each with a newly sampled  $\mathbf{x}_0$ , and the mean membership advantage over the 20 experiments is plotted.

The experiments in Figure 5 are also obtained by approximating the optimal MI adversary with discrete histograms as in the previous paragraph. The only difference is in how the datasets ( $\mathbf{X}$ ,  $\mathbf{y}$ , etc.) are sampled. Specifically, they are sampled according to the distributions for each experiment detailed in Section 5. Again, the histograms are formed by splitting the model’s prediction interval for each  $\gamma$  into 150 equally spaced bins. 20 experiments are performed for each data model, with the means and standard errors reported in the figures.

## E Experimental Verification of Ridge Theory

We verify our theoretical finding that ridge regression increases membership advantage on linear regression models with Gaussian data in the overparameterized regime. The experiment follows the procedure detailed in Section D for Figure 2b except that we only sample 50,000 datasets for each of  $m = 0$  and  $m = 1$  for each  $\gamma$  and each  $\lambda$  for computational efficiency. For this experiment, we set  $n = 100$ ,  $D = 3,000$ , and  $\sigma = 1$ , as in Figure 2b. The results, shown in Figure 7, closely resemble the trend shown in the theoretical plot in Figure 3a, thus verifying our theory.



## References

- [1] N. Carlini, F. Tramèr, E. Wallace, M. Jagielski, A. Herbert-Voss, K. Lee, A. Roberts, T. Brown, D. Song, Ú. Erlingsson, A. Oprea, and C. Raffel, “Extracting training data from large language models,” in *30th USENIX Security Symposium*, pp. 2633–2650, Aug. 2021.
- [2] B. Liu, M. Ding, S. Shaham, W. Rahayu, F. Farokhi, and Z. Lin, “When machine learning meets privacy: A survey and outlook,” *ACM Computing Surveys*, vol. 54, Mar. 2021.
- [3] T. Hastie, A. Montanari, S. Rosset, and R. J. Tibshirani, “Surprises in high-dimensional ridgeless least squares interpolation,” *arXiv preprint arXiv:1903.08560*, 2019.
- [4] W. Fedus, B. Zoph, and N. Shazeer, “Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity,” *arXiv preprint arXiv:2101.03961*, 2021.
- [5] J. Lin, A. Yang, J. Bai, C. Zhou, L. Jiang, X. Jia, A. Wang, J. Zhang, Y. Li, W. Lin, J. Zhou, and H. Yang, “M6-10t: A sharing-delinking paradigm for efficient multi-trillion parameter pretraining,” *arXiv preprint arXiv:2110.03888*, 2021.
- [6] M. Belkin, D. Hsu, S. Ma, and S. Mandal, “Reconciling modern machine-learning practice and the classical bias–variance trade-off,” *Proceedings of the National Academy of Sciences*, vol. 116, no. 32, pp. 15849–15854, 2019.
- [7] Y. Dar, V. Muthukumar, and R. G. Baraniuk, “A farewell to the bias-variance tradeoff? an overview of the theory of overparameterized machine learning,” *arXiv preprint arXiv:2109.02355*, 2021.
- [8] K. Leino and M. Fredrikson, “Stolen memories: Leveraging model memorization for calibrated white-box membership inference,” in *29th USENIX Security Symposium*, pp. 1605–1622, Aug. 2020.
- [9] S. Yeom, I. Giacomelli, M. Fredrikson, and S. Jha, “Privacy risk in machine learning: Analyzing the connection to overfitting,” in *IEEE 31st Computer Security Foundations Symposium*, pp. 268–282, 2018.
- [10] N. Carlini, S. Chien, M. Nasr, S. Song, A. Terzis, and F. Tramer, “Membership inference attacks from first principles,” *arXiv preprint arXiv:2112.03570*, 2021.
- [11] F. Mireshghallah, K. Goyal, A. Uniyal, T. Berg-Kirkpatrick, and R. Shokri, “Quantifying privacy risks of masked language models using membership inference attacks,” *arXiv preprint arXiv:2203.03929*, 2022.
- [12] H. Hu, Z. Salcic, G. Dobbie, and X. Zhang, “Membership inference attacks on machine learning: A survey,” *arXiv preprint arXiv:2103.07853*, 2021.
- [13] R. Shokri, M. Stronati, C. Song, and V. Shmatikov, “Membership inference attacks against machine learning models,” in *IEEE Symposium on Security and Privacy*, pp. 3–18, 2017.
- [14] A. Sablayrolles, M. Douze, C. Schmid, Y. Ollivier, and H. Jégou, “White-box vs black-box: Bayes optimal strategies for membership inference,” in *International Conference on Machine Learning*, pp. 5558–5567, 2019.
- [15] E. Galinkin, “The influence of dropout on membership inference in differentially private models,” *arXiv preprint arXiv:2103.09008*, 2021.
- [16] Y. Wang, C. Wang, Z. Wang, S. Zhou, H. Liu, J. Bi, C. Ding, and S. Rajasekaran, “Against membership inference attack: Pruning is all you need,” *arXiv preprint arXiv:2008.13578*, 2020.
- [17] Y. Long, V. Bindschaedler, and C. A. Gunter, “Towards measuring membership privacy,” *arXiv preprint arXiv:1712.09136*, 2017.

- [18] A. Salem, Y. Zhang, M. Humbert, P. Berrang, M. Fritz, and M. Backes, “MI-leaks: Model and data independent membership inference attacks and defenses on machine learning models,” *arXiv preprint arXiv:1806.01246*, 2018.
- [19] C. A. Choquette-Choo, F. Tramèr, N. Carlini, and N. Papernot, “Label-only membership inference attacks,” in *Proceedings of the 38th International Conference on Machine Learning*, vol. 139, pp. 1964–1974, July 2021.
- [20] Y. Kaya and T. Dumitras, “When does data augmentation help with membership inference attacks?,” in *International Conference on Machine Learning*, pp. 5345–5355, 2021.
- [21] A. Bagmar, S. R. Maiya, S. Bidwalka, and A. Deshpande, “Membership inference attacks on lottery ticket networks,” *arXiv preprint arXiv:2108.03506*, 2021.
- [22] S. Rezaei, Z. Shafiq, and X. Liu, “Accuracy-privacy trade-off in deep ensemble,” *arXiv preprint arXiv:2105.05381*, 2021.
- [23] M. Belkin, D. Hsu, and J. Xu, “Two models of double descent for weak features,” *SIAM Journal on Mathematics of Data Science*, vol. 2, no. 4, pp. 1167–1180, 2020.
- [24] C. Dwork, “Differential privacy: A survey of results,” in *International Conference on Theory and Applications of Models of Computation*, pp. 1–19, 2008.
- [25] M. Bun, J. Ullman, and S. Vadhan, “Fingerprinting codes and the price of approximate differential privacy,” *SIAM Journal on Computing*, vol. 47, no. 5, pp. 1888–1938, 2018.
- [26] P. L. Bartlett, P. M. Long, G. Lugosi, and A. Tsigler, “Benign overfitting in linear regression,” *Proceedings of the National Academy of Sciences*, vol. 117, no. 48, pp. 30063–30070, 2020.
- [27] A. Tsigler and P. L. Bartlett, “Benign overfitting in ridge regression,” *arXiv preprint arXiv:2009.14286*, 2020.
- [28] B. Hui, Y. Yang, H. Yuan, P. Burlina, N. Z. Gong, and Y. Cao, “Practical blind membership inference attack via differential comparisons,” *arXiv preprint arXiv:2101.01341*, 2021.
- [29] J. Ye, A. Maddi, S. K. Murakonda, and R. Shokri, “Enhanced membership inference attacks against machine learning models,” *arXiv preprint arXiv:2111.09679*, 2021.
- [30] Y. Kaya, S. Hong, and T. Dumitras, “On the effectiveness of regularization against membership inference attacks,” *arXiv preprint arXiv:2006.05336*, 2020.
- [31] R. Sarathy and K. Muralidhar, “Evaluating Laplace noise addition to satisfy differential privacy for numeric data,” *Transactions on Data Privacy*, vol. 4, no. 1, pp. 1–17, 2011.
- [32] A. Rahimi and B. Recht, “Random features for large-scale kernel machines,” in *Proceedings of the 20th International Conference on Neural Information Processing Systems*, pp. 1177–1184, 2007.
- [33] Y. Cho, *Kernel Methods for Deep Learning*. PhD thesis, University of California, San Diego, 2012.
- [34] F. Rubio and X. Mestre, “Spectral convergence for a general class of random matrices,” *Statistics & Probability Letters*, vol. 81, no. 5, pp. 592–602, 2011.
- [35] E. Dobriban and Y. Sheng, “Wonder: Weighted one-shot distributed ridge regression in high dimensions,” *Journal of Machine Learning Research*, vol. 21, pp. 66–1, 2020.
- [36] E. Dobriban and Y. Sheng, “Distributed linear regression by averaging,” *The Annals of Statistics*, vol. 49, no. 2, pp. 918–943, 2021.
- [37] R. D. Cook and L. Forzani, “On the mean and variance of the generalized inverse of a singular Wishart matrix,” *Electronic Journal of Statistics*, vol. 5, pp. 146–158, 2011.